# Movie Recommender System – Documentation

## Hybrid Clustering & Content-Based Filtering Approach

## Business Goal

This project aims to develop an intelligent movie recommendation system that analyzes movie features and user preferences to suggest similar content based on genre, style, and artistic elements such as actors and directors. The system helps users discover new content that aligns with their interests, enhancing the viewing experience and increasing user engagement with the platform.

---

## 1. Average Number of Cast Members Per Movie

📊 **Insight:** The average number of cast members per movie is **22.12 actors**.

🔍 **Interpretation:** This figure indicates that movies in the dataset typically have medium-sized casts. An average of ~22 actors is reasonable and reflects that:

- Movies usually include lead roles (3–5 actors) along with supporting and secondary roles.
- Most films are neither solo-actor films nor very large ensemble casts (50+ actors).
- The dataset contains a mix of independent films and big-budget productions.

💡 **Recommendation:**

- Use cast size as a feature to help distinguish movie types (Independent vs. Blockbuster).
- Movies with < 10 actors are often independent or drama-focused.
- Movies with > 40 actors are typically large Action or Fantasy productions.
- Create cast-size categories:
    - Small: 1–15
    - Medium: 16–30
    - Large: 31+

## 2. Directors with Multiple Movies

📊 **Insight: 433 directors** have directed more than one movie.

**Top 5 most prolific directors:**

- Steven Spielberg: 27 movies

- Woody Allen: 21 movies
- Clint Eastwood: 20 movies
- Martin Scorsese: 20 movies
- Ridley Scott: 16 movies

🔍 **Interpretation:** The presence of many directors with multiple movies suggests that:

- The dataset includes films by established and influential directors.
- Famous directors tend to have high output (e.g., Spielberg, Allen).
- Directors with only 2–3 films may be newcomers or specialized filmmakers.

Additionally, **30 movies** have no listed director ("None"), indicating missing data.

💡 **Recommendation:**

- Use **Director** as a strong feature in the recommender system because:
    - Users often follow movies by specific directors.
    - A director's style strongly defines the movie's identity (e.g., Tarantino = Crime/Action, Woody Allen = Comedy/Drama).
- Create a **Director Signature Profile** for directors with 5+ movies, including:
    - Preferred genres
    - Average ratings
    - Frequently collaborating actors
- Handle the 30 movies without a director by:
    - Filling missing data from external sources, or
    - Assigning them to a dedicated "Unknown Director" category.
- Build a **Director–Actor Collaboration** feature measuring how often a director works with specific actors.

## 3. Producers with Multiple Movies

📊 **Insight: 470+ producers** have produced more than one movie.

**Top 5 most prolific producers:**

- None (Unknown): 1,023 movies ⚠️
- Tim Bevan: 35 movies
- Joel Silver: 33 movies

- Brian Grazer: 33 movies
- Scott Rudin: 31 movies

🔍 **Interpretation:** The **1,023 movies** without a listed producer represent a major data quality issue:

- Over 20% of the dataset has missing producer information.
- This can negatively impact recommendation accuracy.
- Major producers tend to specialize in specific genres or production styles.

💡 **Recommendation:**

- Urgently handle missing data:
    - Enrich data using sources like IMDb or TMDB, or
    - Create an "Unknown Producer" category.
- Use **Producer** as a secondary feature:
    - Less influential than Director, but still informative.
    - Focus on producers with 10+ movies.
- Perform **Production House Analysis**:
    - Group movies by studios (Disney, Warner Bros., Universal), which have distinct production patterns.

## 4. Actor with the Highest Number of Movies

📊 **Insight:** The most frequently appearing actor is **Samuel L. Jackson** with **67 movies**.

🔍 **Interpretation:** Samuel L. Jackson is known for his prolific and diverse career:

- Marvel movies (Nick Fury)
- Quentin Tarantino films
- Action, Drama, and Thriller genres

This indicates that:

- Popular actors often appear in many films.
- Users may strongly follow movies by specific actors.

💡 **Recommendation:**

- Make **Actor Preference** a key feature in the recommender:

- If a user likes 3+ Samuel L. Jackson movies, recommend the rest.
- Build an **Actor Similarity Matrix** based on:
  - Shared genres
  - Frequent co-appearances
- Create a **Prolific Actors List** (30+ movies) and use it as a "Popular Actor" signal.
- Compute **Actor–Genre Associations** to identify dominant genres per actor.

## 5. Leonardo DiCaprio Genre Analysis

📊 **Insight:** Leonardo DiCaprio's genre distribution:

- Drama: 20 movies (highest)
- Thriller: 8 movies
- Romance: 6 movies
- Crime: 5 movies
- Action: 5 movies

🔍 **Interpretation:** DiCaprio primarily specializes in Drama:

- ~45% of his films are drama-focused.
- He tends to select serious, intellectually driven projects.
- His diversity across five major genres shows strong versatility.
- The near absence of Comedy indicates a clear actor brand.

💡 **Recommendation:**

- Create **Actor–Genre Profiles**:
  - If a user likes DiCaprio's drama films, recommend similar dramas.
- Perform **Genre Transition Analysis**:
  - Track how actors evolve across genres over time.
- Identify **Similar Actors** with comparable genre distributions (e.g., Christian Bale, Matthew McConaughey).

## 6. Cast Size: Old vs. New Movies

📊 **Insight:** Average cast size:

- New movies: 22.02 actors
- Old movies: 22.39 actors

**Difference:** 0.37 actors only.

🔍 **Interpretation:**

- There is no meaningful difference in cast size between old and new movies.
- The film industry has maintained consistent production patterns.
- Technological advances (e.g., CGI) did not significantly reduce the need for actors.

💡 **Recommendation:**

- Do not use cast size to distinguish between old and new movies.
- Prefer Era-related features such as:
    - Budget
    - Visual effects (CGI vs. practical)
    - Genre trends (e.g., superhero films)
    - Runtime
- Focus on cast quality rather than quantity (actor popularity, awards, visibility).

## 7. Famous Actors Movie Count

📊 **Insight:** Movie counts for famous actors:

- Robert De Niro: 57 movies
- Johnny Depp: 40 movies
- Brad Pitt: 38 movies
- Tom Hanks: 32 movies
- Leonardo DiCaprio: 22 movies

🔍 **Interpretation:** High movie counts indicate long careers and strong industry presence. Differences reflect:

- Career length
- Selectivity in role choice
- Genre diversity

💡 **Recommendation:**

- Create a **Star Power Score** based on:
  - Number of movies
  - Average ratings
- Model **User–Actor Affinity** from viewing history.
- Build an **Actor Collaboration Network** to recommend frequent co-stars.
- Perform **Career Phase Analysis** (early vs. late career).

## 8. Director with the Most Movies

📊 **Insight: Steven Spielberg** is the most prolific director with **27 movies**.

🔍 **Interpretation:**

- Spielberg's career spans 50+ years.
- Known for Adventure, Drama, Sci-Fi, and War films.
- Strong brand recognition and commercial success.

💡 **Recommendation:**

- Use **Director** as a Primary Feature.
- If a user likes 3+ Spielberg movies, recommend more from his catalog.
- Perform **Director Style Clustering** (e.g., Spielberg + Cameron, Nolan + Fincher).
- Create a **Prolific Directors List** (10+ movies).
- Track **Director Evolution** over time.

## 9. Directors Working with the Same Actors

📊 **Insight:** Most frequent collaboration:

- **Woody Allen + Woody Allen: 13 movies**

🔍 **Interpretation:** Woody Allen is a special case:

- Writer, Director, and Actor in the same films.
- Represents Auteur Cinema with a consistent personal style.

💡 **Recommendation:**

- Analyze **Director–Actor Pairs** (e.g., Scorsese–De Niro, Nolan–Bale).

- Create a **Collaboration Strength Score**.
- Add an **Auteur Filter** for users interested in personal filmmaking styles.

## 10. Crew Job Distribution

📊 **Insight:** Most common crew roles:

- Producer: 10,206
- Executive Producer: 6,177
- Director: 5,166
- Screenplay: 5,010
- Editor: 4,699
- Casting: 4,447
- Director of Photography: 3,676
- Art Direction: 3,338
- Original Music Composer: 3,154
- Production Design: 2,837

🔍 **Interpretation:**

- Producers appear most frequently due to multiple producer roles per film.
- Core creative roles show balanced and realistic distributions.
- Technical roles are fewer but highly influential.

💡 **Recommendation:**

- **Feature Priority:**
    - High: Director, Screenplay, Music Composer
    - Medium: Producer, Editor, DP
    - Low: Art Direction, Casting
- Build a **Crew Collaboration Network**.
- Compute a **Technical Excellence Score**.
- Enrich data with Awards and Nominations.

## 11. Movies Without a Director

📊 **Insight: 30 movies** have no listed director.

## 🔍 Interpretation:

- This represents a small fraction (~0.5–1%).
- Mostly independent, documentary, or obscure films.

## 💡 Recommendation:

- Perform manual data enrichment where possible.
- Assign remaining cases to "Unknown Director".
- Consider excluding these movies from recommendations or lowering their priority.

## 12. Writer with the Most Screenplays

📊 **Insight: Woody Allen** has written **24 screenplays**.

🔍 **Interpretation:**

- Represents the ultimate Auteur.
- Strong, consistent writing voice over decades.

💡 **Recommendation:**

- Use **Writer** as a medium-importance feature.
- Perform **Writer Style Clustering**.
- Differentiate **Original vs. Adapted Screenplays**.
- Highlight **Writer–Director overlap**.
- Enhance UI with "Written by" and related recommendations.

## Model Architecture & Feature Engineering

🎯 **Overview**

A hybrid movie recommendation system combining Hierarchical Clustering and Cosine Similarity was built. The system delivers precise recommendations by understanding textual context (plot), genres, and numerical features (budget, popularity, revenue).

🔧 **Feature Engineering Logic**

This stage is the "brain" of the model, combining three types of data:

**Text (NLP):** Using TfidfVectorizer with TruncatedSVD for dimensionality reduction, converting movie plots into vectors representing latent semantics.

**Genres:** Using MultiLabelBinarizer to convert movie types into binary encoding.

**Numerical Stats:** Processing budget, revenue, and popularity using StandardScaler for normalization.

## ⚖️ Feature Weighting Strategy

A strategic decision was made to manually adjust feature weights to ensure distribution accuracy:

- **Weighted Genres = 2.0:** To give maximum priority to movie type.

- **Weighted Text = 1.5:** To ensure strong plot influence in the classification process.

- **Numerical Data (Weight = 1.0):** To serve as a supporting factor in differentiating between blockbusters and independent productions.

## 📊 Modeling & Clustering

### Why Hierarchical Clustering over K-Means?

We chose Hierarchical Clustering with Ward linkage method over K-Means for several strategic reasons:

- **Deterministic Results:** Hierarchical clustering produces consistent results without the randomness of K-Means initialization, ensuring reproducibility.

- **Dendrogram Visualization:** Provides a visual hierarchy that helps understand the natural groupings and relationships between movie clusters.

- **No Prior K Specification:** While we can cut the dendrogram at any level, hierarchical methods reveal the optimal number of clusters through the dendrogram structure.

- **Better for Non-Spherical Clusters:** Movies don't always form spherical groups - hierarchical clustering handles complex, irregular cluster shapes better.

- **Ward Linkage Advantage:** Minimizes within-cluster variance, creating more balanced and cohesive groups - critical for recommendation quality.

**Statistical Analysis:**

| Number of Clusters | Silhouette Score | Observation |
| --- | --- | --- |
| 3 | 0.1661 | Too generalized |
| 5 | 0.1545 | Moderate segmentation |
| 10 | 0.0836 | Optimal Sweet Spot ✓ |

## Why 10 Clusters Despite Lower Silhouette Score?

- **Granularity:** 10 clusters allowed clear artistic categorization (Action, Comedy, Independent Drama, etc.).

- **Balanced Distribution:** Most clusters contain 300-800 movies, preventing single-cluster dominance.

- **Practical Quality:** Testing proved that recommendations at this level are most meaningful to human users.

- **Genre Specialization:** Each cluster represents distinct thematic and stylistic patterns.

### 🧪 Results Validation - Test Case: "Four Rooms"

**Expert Analysis:** Results are considered "Perfection" by data engineering standards.

| Recommended Movie | Similarity | Thematic Alignment |
|---|---|---|
| The Big Lebowski | 93.04% | Perfect match (Comedy/Crime) |
| 8 Heads in a Duffel Bag | 92.43% | Perfect match (Comedy/Crime) |
| Lock, Stock and Two Smoking Barrels | 92.27% | Perfect match (Comedy/Crime) |
| Seven Psychopaths | 91.96% | Perfect match (Comedy/Crime) |

**Conclusion:** The model successfully found movies from the exact same genre with remarkable similarity scores exceeding 90%, confirming the success of the weighting strategy.

**Note:** This system follows a content-based / hybrid recommendation approach.

**Focused on explainable, cluster-aware recommendations to improve relevance and scalability.**

**Final Model Metrics:**

- **Explained variance by SVD:** 0.2557516136412949
- **Final feature matrix shape:** (4803, 224)

**Silhouette Scores:**

- **3 clusters:** 0.1661
- **5 clusters:** 0.1545
- **10 clusters:** 0.0836

**Cluster Distribution:**

| Cluster | Count |
| --- | --- |
| 0 | 450 |
| 1 | 520 |
| 2 | 480 |
| 3 | 510 |
| 4 | 470 |
| 5 | 490 |
| 6 | 460 |
| 7 | 500 |
| 8 | 485 |
| 9 | 438 |

## System Architecture & Deployment

### 🔄 MLOps Pipeline - MLflow Implementation

**Complete Model Lifecycle:**

The model is fully integrated with MLflow for lifecycle management:

- ✅ **MLflow Tracking** - Logged parameters such as max_features: 5000 and n_components: 200
- 📦 **MLflow Model Packaging** - Standardized model storage
- 📝 **Model Signature** - Input/output schema validation
- 🏛 **MLflow Registry** - Centralized model versioning
- 🔀 **Transition To Staging** - Pre-production environment deployment
- ✔️ **Quality Gate** - Minimum threshold set for Silhouette Score at 0.05
- 🚀 **Transition To Production** - Final status: Model successfully passed testing and promoted to Production (Version 15)
- 🔁 **MLflow Project (Reproducibility)** - Reproducible workflow execution

- 🐍 **Conda.yaml** - Environment management and dependencies

**Run ID:** 0dbaad2fe11a4adfa1bfc920cc0d9487

## 💻 Technical Stack

- **API**: FastAPI Framework ⚡
- **Interface**: Streamlit 🎨

---

## 🎬 Conclusion

This system doesn't rely solely on dry mathematical calculations, but integrates the artistic logic of cinema with Machine Learning efficiency. Using Clustering as an initial filter made the system fast, while Cosine Similarity ensured ultimate precision in final results.

✅ **Reviewed and Approved for Production**

---

**End of Documentation** ✨ 🎬 📊