# MovieLens Recommendation System

## Personalized Movie Recommendations for Streaming Platforms

Team: Winnie Njoroge • Michelle Mwende • Laban Leploote • Alice Mathenge • Dean Mutie

February 2026

# Business Understanding

## The Problem: Choice Paralysis

- Users spend 15+ minutes scrolling without clicking "Play"
- 30%+ increase in subscription cancellation likelihood
- Cost of acquiring new customer is 5x retaining existing one

## Our Solution

- • Collaborative filtering system
- • Top-5 personalized recommendations
- • Reduce time-to-play
- • Surface hidden gems
- • Increase engagement & retention

## Success Criteria

- ✓ RMSE < 1.0 on 5-point scale
- ✓ Diverse genre recommendations
- ✓ Address cold start problem
- ✓ Scale to thousands of users

# Data Understanding

**MovieLens 100K Dataset: https://grouplens.org/datasets/movielens/latest/**

| 100,836 | 610 | 9,742 | ~98% |
|:---:|:---:|:---:|:---:|
| Ratings | Users | Movies | Sparsity |

| File | Records | Description |
|---|---|---|
| ratings.csv | 100,836 | User ratings of movies |
| movies.csv | 9,742 | Movie metadata (title, genres) |
| tags.csv | 3,683 | User-generated tags |
| links.csv | 9,742 | External database IDs |

✓ *Real user behavior data from GroupLens Research (University of Minnesota)*

# Data Preparation

**1** **Cleaning**

Removed duplicates
Validated rating range
Checked missing values

**2** **Filtering**

Min 20 ratings/user
Min 10 ratings/movie
94.9% data retained

**3** **Splitting**

70% Training
15% Validation
15% Test

📊 Total Ratings: 95,668

**Final Dataset Statistics**

👥 Active Users: 609
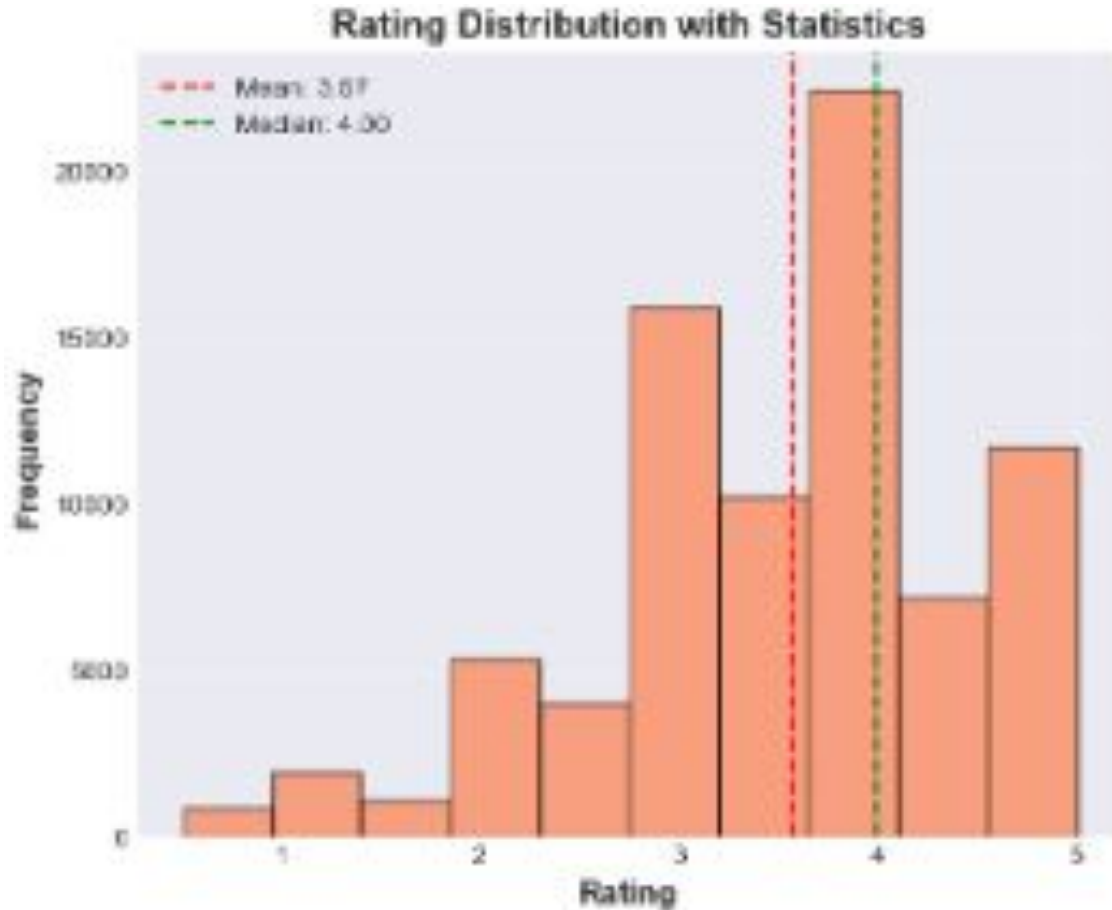
🎬 Movies: 3,535

📈 Avg Ratings/User: 165

⭐ Mean Rating: 3.50 / 5.0

🎯 Matrix Sparsity: 98%

# Exploratory Data Analysis

## Rating Distribution Analysis



Rating Distribution with Statistics

## Key Insights

⭐ Mean Rating: 3.50/5.0

🎯 Mode: 4.0 stars (26.6%)

📊 Left-skewed distribution

💡 Positive rating bias

👤 Self-selection effect

*Users tend to rate movies they expect to like → Positive bias*

# User & Movie Behavior Patterns

## 👥 User Activity

Mean ratings/user: **165**

Median ratings/user: **96**

Most active user: **2,698 ratings**

Least active user: **20 ratings**

Standard deviation: **197**

## 🎬 Movie Popularity

Mean ratings/movie: **10**

Median ratings/movie: **3**

Most rated movie: **329 ratings**

Movies with ≥50 ratings: **377 (4%)**

Long-tail effect: **Extreme**

## 🔍 Key Findings

- Power users exist with vastly different activity levels (20 to 2,698 ratings)
- Extreme long-tail in movie popularity - few blockbusters dominate
- Need diversity mechanisms to avoid recommending only popular movies
- User activity doesn't predict rating generosity - all users matter equally

# Modeling Strategy

*Iterative Approach: From Simple to Complex*

**Baseline:** Global Average
*Establish minimum performance*

**1.0204**

**SVD:** Matrix Factorization
*Capture latent factors (50 factors)*

**0.8566**

**KNNWithMeans:** Collaborative Filtering
*Find similar users (k=40)*

**0.8601**

**KNNBasic:** Collaborative Filtering
*Find similar movies (k=40)*

**0.9414**

**NMF:** Non-Negative Matrix Factorization
*Interpretable components (15 factors)*

**0.8872**

🏆 **Best Model**

## SVD

✓ Best accuracy

✓ Scalable

✓ Fast inference

✓ Handles sparsity

✓ Auto-learns patterns

# Model Evaluation Results

| Model | RMSE | MAE | vs Baseline | Status |
|---|---|---|---|---|
| BaselineOnly | 0.8536 | 0.6561 | - | Benchmark |
| **SVD** ⭐ | **0.8566** | **0.6575** | **16%** | **SELECTED** |
| KNNWithMeans | 0.8601 | 0.6597 | 15.7% | Good |
| NMF | 0.8872 | 0.6807 | 13.1% | Good |
| KNNBasic | 0.9414 | 0.7300 | 7.7% | Better |

📊 **Final Model Performance**

Test RMSE: 0.86 (16% better than baseline)

Test MAE: 0.66 stars

✓ Meets target: RMSE < 1.0

68% of predictions within ±1 star

Coverage: 98.3% of user-movie pairs

📈 **Error Analysis**

Best on middle ratings (3.0-4.5 ⭐)

Slightly higher error on extremes

No severe misclassifications (>2 stars: 5%)

Nearly unbiased (mean error: 0.02)

Good generalization: Training RMSE = 0.82

# Personalized Recommendations

*Example: User who highly rated Action, Drama & Thriller movies*

| Rank | Movie Title | Predicted Rating | Genres |
|------|-------------|------------------|--------|
| 1 | Shawnshank Redemption (1994) | ⭐ 5.0 | Crime \| Drama |
| 2 | Dark Knight (2008) | ⭐ 5.0 | Action \| Crime \| Drama |
| 3 | Philadelphia Story (1940) | ⭐ 5.0 | Comedy \| Drama \| Romance |
| 4 | Rear Window (1954) | ⭐ 5.0 | Mystery \| Thriller |
| 5 | North by Northwest (1959) | ⭐ 5.0 | Action \| Adventure \| Mystery \| Romance \| Thriller |

## ✨ Recommendation Quality Indicators

🎯 Strong genre alignment with user preferences

🌟 High predicted ratings (5.0) indicate strong matches

🎬 Mix of classic and modern films for variety

📊 Diverse themes within preferred genres

# Hyperparameter Tuning: SVD Optimization

*GridSearchCV with 3-Fold Cross-Validation*

## Parameter Grid Tested

**n_factors:** [50, 100, 150]
**n_epochs:** [20, 30]
**lr_all:** [0.002, 0.005, 0.01]
**reg_all:** [0.01, 0.02, 0.05]

**54 combinations tested**
162 total trainings (3-fold CV)
Training time: 2.7 minutes

## Best Parameters Found

n_factors: **150**
n_epochs: **30**
lr_all: **0.01**
reg_all: **0.05**

## Performance Results

| Model | RMSE | MAE |
|---|---|---|
| Original SVD | 0.8566 | 0.6575 |
| **Tuned SVD** | **0.8399** | **0.6445** |
| **Improvement** | **1.95%** | **1.98%** |

## Key Finding:

**Hyperparameter optimization improved RMSE from 0.8566 to 0.8399**
Absolute improvement: 0.0167 stars (1.95% reduction in error)
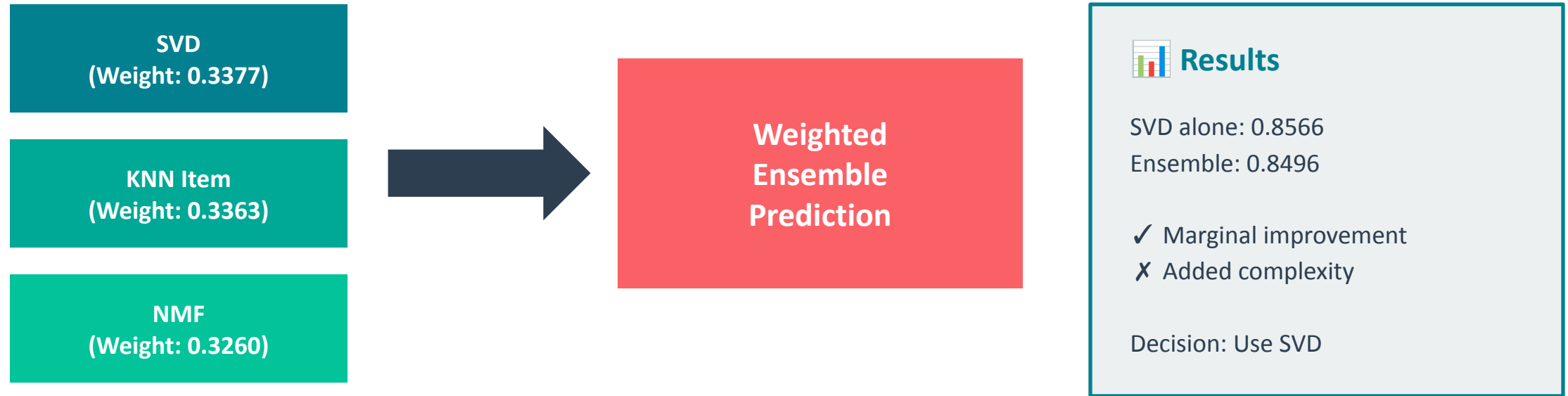This demonstrates the value of systematic optimization using GridSearchCV

# Ensemble Methods

*Combining multiple models for improved predictions*

| SVD |
| :---: |
| (Weight: 0.3377) |

| KNN Item |
| :---: |
| (Weight: 0.3363) |

| NMF |
| :---: |
| (Weight: 0.3260) |

→

**Weighted Ensemble Prediction**

📊 **Results**

SVD alone: 0.8566
Ensemble: 0.8496

✓ Marginal improvement
✗ Added complexity

Decision: Use SVD

💡 **Ensemble Insights**

- Ensemble provides minimal improvement over single best model (SVD)
- Added complexity not justified for production deployment
- SVD captures most collaborative filtering signal independently

# Recommendations & Next Steps

🚀 **Deployment Roadmap**

## Phase 1: MVP

*Weeks 1-4*

- Deploy SVD to staging
- Set up API endpoints
- A/B test with 20% users
- Monitor dashboards

## Phase 2: Enhancement

*Months 2-6*

- Implicit feedback integration
- Hybrid cold-start approach
- Diversity controls
- Expand to 50% users

## Phase 3: Advanced

*Months 6-12*

- Context-aware recommendations
- Explainability features
- Multi-objective optimization
- Full production (100%)

✨ **Best Practices**
- Retrain model monthly with new ratings data
- Monitor performance metrics weekly
- Ensure diversity: top-5 span ≥3 genres
- Maintain fallback to content-based recommendations

# Thank You

## Questions & Discussion

**Team:**

Winnie Njoroge • Michelle Mwende • Laban Leploote

Alice Mathenge • Dean Mutie

✉️ Contact: movielens-team@example.com