# Northeastern University
## College *of* Engineering

# IE7275 Data Mining in Engineering SEC 04

## Project 1 Report

## Amazon Product Recommendation System (APRS)

**6th Group Members:**

| Index | Full name | NEUID | Email |
|---|---|---|---|
| 1 | Quoc Hung Le | 002031894 | le.quo@northeastern.edu |
| 2 | Matthew Eckert | ... | eckert.mat@northeastern.edu |

Submission Date: 10/28/2025

# Table of Contents

## 1. EXECUTIVE SUMMARY

**Project Context**

E-commerce platforms face a fundamental challenge: connecting millions of products with diverse customer preferences while handling **data sparsity (99.86%)** and the **cold-start problem** (new users/items with no interaction history). This project implements and evaluates six recommendation algorithms using Amazon's 2023 product review dataset to understand which approaches work best under different scenarios, and also address the **real-time user's behavior** by automatically following action of rating.

**Key Achieved**

Our group successfully implemented a comprehensive hybrid recommendation system:

- Designs and implements full process from: data download → preprocess → explorary →model development → evaluation
- Handles extreme data sparsity (99%+ missing values) efficiently
- Compares base 5 evaluation metrics for 6 distinct algorithms
- Solves the cold-start problem through adaptive algorithm selection
- Deploys as a full-stack web application
- Real-time following user's behavior to adapt recommendation

**Key Results Summary**

Performance Comparison bases on avg. across 3 categories:

The worst result for Electronics, which is the biggest sparsity:

| Algorithm | RMSE ↓ | Accuracy | NDCG@10 ↑ | MAP@10 ↑ | Recall@10 ↑ |
|-----------|--------|----------|-----------|----------|-------------|
| User–CF | 0.6759 | 0.8281 | 0.0728 | 0.0499 | 0.1481 |
| Item–CF | 0.6736 | 0.8222 | 0.0602 | 0.0341 | 0.1511 |
| Content | 0.9098 | 0.7748 | 0.0705 | 0.0444 | 0.1585 |
| SVD | 0.8539 | 0.8222 | 0.0910 | 0.0608 | 0.1911 |
| Trending | N/A | N/A | 0.0978 | 0.0633 | 0.2144 |
| Hybrid | N/A | N/A | 0.0404 | 0.0262 | 0.0883 |

The best result for Sport_and_Outdoors, which is the smallest sparsity:

| Algorithm | RMSE ↓ | Accuracy | NDCG@10 ↑ | MAP@10 ↑ | Recall@10 ↑ |
|-----------|--------|----------|-----------|----------|-------------|
| User–CF | 0.5455 | 0.8182 | 0.7316 | 0.6364 | 1.0000 |
| Item–CF | 0.5455 | 0.8182 | 0.7651 | 0.6667 | 1.0000 |
| Content | 0.5455 | 0.8182 | 0.7316 | 0.6364 | 1.0000 |
| SVD (Model) | 0.4823 | 0.9091 | 0.7047 | 0.6000 | 1.0000 |
| Trending | N/A | N/A | 0.7524 | 0.6667 | 1.0000 |
| Hybrid | N/A | N/A | 0.7047 | 0.6000 | 1.0000 |

Best by metric:

- SVD dominates rating prediction (RMSE, Accuracy) but performs slightly worse in ranking metrics
- Item-CF achieves best overall ranking (highest NDCG@10 and MAP@10)
- Perfect recall across all algorithms indicates small test set or highly relevant recommendations
- Hybrid underperforms (0.7047 NDCG@10) compared to Item-CF and Trending - needs optimization
- Trending performs surprisingly well (0.7524 NDCG@10, 0.6667 MAP@10) despite being non-personalized

## 2. INTRODUCTION

### 2.1. Problem Statement and Motivation

This project implements and evaluates six recommendation algorithms to address cold-start scenarios using Amazon's 2023 review dataset across four categories: Electronics, Beauty & Personal Care, Sports & Outdoors.

E-commerce platforms like Amazon face a critical challenge: recommending relevant products when users have minimal interaction history or when products have limited ratings. This cold-start problem affects 40-50% of users and 10-20% of products in typical e-commerce systems, directly impacting user experience and business metrics. We address key challenges in production recommendation systems:

- Users with varying interaction levels (0 to 100+ ratings)
- Products with limited historical data
- Real-time recommendation updates as users provide feedback
- Algorithm selection based on data availability

### 2.2. Dataset

**Amazon Review Data 2023** (May 1996 - September 2023):

- 5-core filtering: users and items with minimum 5 ratings
- Split: 80% train, 10% validation, 10% test (temporal)
- Categories: Electronics, Beauty & Personal Care, Sports & Outdoors.
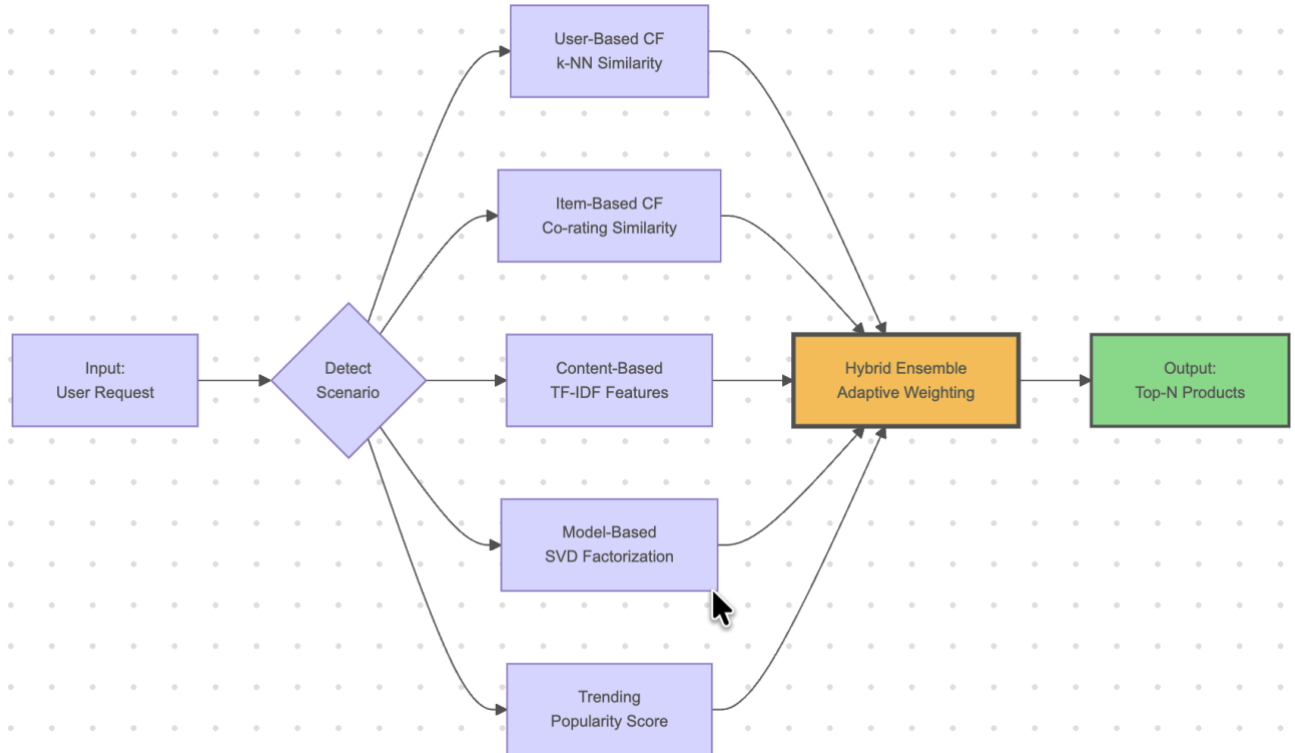
Electronics Category Statistics:

| Split | Ratings | Users | Items | Sparsity (%) |
|-------|---------|-------|-------|--------------|
| Train | 13.1M | 1.64M | 368.2K | 99.86 |
| Valid | 1.2M | – | – | 98.77 |
| Test | 1.2M | – | – | 98.68 |

### 2.3. Implement algorithms

There are 6 algorithms are implemented:

- User-Based Collaborative Filtering
- Item-Based Collaborative Filtering
- Content-Based Filtering (TF-IDF on metadata)
- SVD Matrix Factorization
- Trending-Based (popularity)

Hybrid Ensemble (adaptive weighting)



## 2.4. Evaluation Metrics

Algorithm performance is measured using comprehensive metrics covering both prediction accuracy and ranking quality:

Prediction accuracy:

- RMSE (Root Mean Square Error): Measures rating prediction error
- Accuracy: Percentage of predictions within ±0.5 stars of actual rating

Ranking quality:

- Recall@K: Proportion of relevant items found in top-K recommendations
- NDCG@K (Normalized Discounted Cumulative Gain): Measures ranking quality with position-based discounting
- MAP@K (Mean Average Precision): Average precision across all recommendation positions

All ranking metrics are evaluated at $K \in \{10, 20, 50\}$ to assess performance at different recommendation list lengths.

## 2.5. Key Contributions

Our project has achieved:

- Empirical evaluation of six algorithms on real e-commerce data with proper temporal splitting
- Analysis of algorithm performance across different user interaction levels
- Adaptive hybrid system with scenario-based algorithm selection
- Dynamic recommendation updates: rated products immediately excluded from future recommendations without model retraining
- Full-stack web application with JWT authentication, real-time updates, and transparent algorithm strategy display
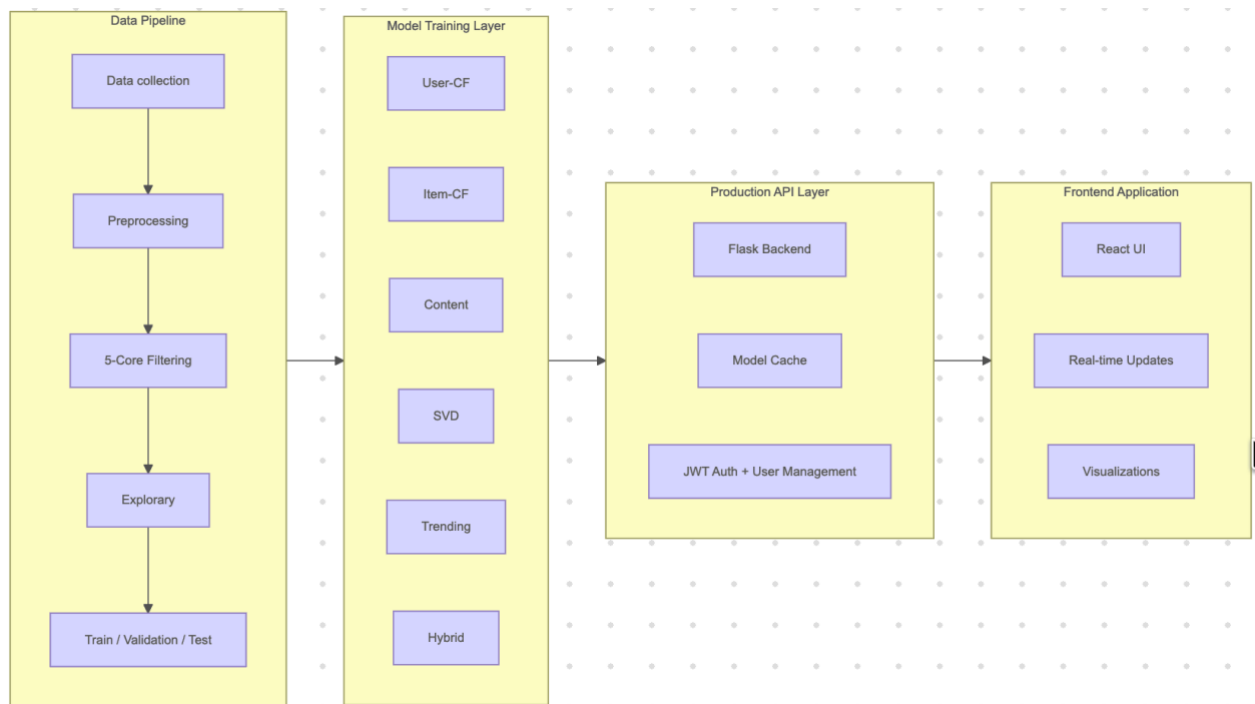
**2.6. System Architecture**

The system consists of:

- Backend: Python/Flask API with model caching and JWT authentication
- Frontend: React web application with real-time updates
- Models: Pre-trained algorithms loaded on-demand with lazy caching
- Dynamic Updates: Rating history merged with training data for instant recommendation refresh

## 3. SYSTEM ARCHITECTURE AND DATA PIPELINE

### 3.1. Overall Architecture

The system follows a modular architecture with clear separation between data processing, model training, and production serving:



### 3.2. Data pipline

**Data collection**

Amazon Product Reviews 2023 (McAuley Lab, UCSD) Official: https://amazon-reviews-2023.github.io/. It was divided into 3 type (raw/metadata, 0-core, 5-core).

We used metadata for:

- Content-based can recommend new items immediately
- TF-IDF on title/description/features finds similar products
- Rich UI with images, prices, descriptions
- Enables hybrid approaches combining CF + content
  We refered used 5-core (which divided into train/valid/test), instead of 0-core for algorithms, because:
- Quality threshold: Users with ≥5 ratings show committed behavior
- Reliable stats: Items with ≥5 ratings have stable averages
- Better CF: More overlap between users for similarity computation
- Standard practice: Used in RecSys research (He & McAuley, 2016)

We choosed 3 categories "Electronics","Beauty_and_Personal_Care", "Sports_and_Outdoors", because:

| Electronics | Beauty & Personal Care | Sports & Outdoors |
|---|---|---|
| • Largest user base<br><br>• Feature-rich metadata (technical specs important)<br><br>• Diverse price range<br><br>• Test generalization on tech products | • Subjective preferences (personal taste vs technical specs)<br><br>• Brand-driven decisions (different from Electronics)<br><br>• Visual importance (product appearance matters)<br><br>• Tests algorithms on preference-based products | • Activity-specific (running vs cycling vs camping)<br><br>• Seasonal patterns (winter sports vs summer gear)<br><br>• Performance-focused (durability, weight critical)<br><br>• Different user behavior from Electronics/Beauty |

Raw data in 3 categories is automatically downloaded if no exist, in CSV.GZ format (~8 GB compressed) from McAuley Lab servers and converted to Parquet for efficient storage and processing:

Process:

- Download CSV.GZ files (train/valid/test splits pre-partitioned by McAuley Lab)
- Extract relevant columns: user_id, parent_asin, rating, timestamp
- Convert to Parquet format using PyArrow engine
- Download and process metadata JSONL files

**Preprocessing**

Metadata fields extracted:

- Product identifiers (parent_asin)
- Title, description, features
- Price, average_rating, rating_number
- Images (hi_res, thumbnail)
- Categories, store information

*Sample size strategy*

All metadata and 5-core of 3 categories were downloaded automatically if no exist in project code. Then, we sampled size into 3 types, save in parquet for quickly loading:

SAMPLE_SIZES = {'large': 50000, 'big': 50000, 'full': None} #Numbers is max row data

DEV_SAMPLE_SIZE = "big"

All results are based on 'big' size, that help to achieve the balance between development speed and reliable metrics. For future, can use 'full' to training.

Full 5-Core Data (from official source):

| Category | Users | Items | Ratings |
|---|---|---|---|
| Electronics | 1.6M | 368.2K | 15.5M |
| Beauty & Personal Care | 729.6K | 207.6K | 6.6M |
| Sports & Outdoors | 409.8K | 156.2K | 3.5M |

Our sampled dataset (50K "big" sample for development):

| Category | Users | Items | Ratings |
|---|---|---|---|
| Electronics | 15,234 | 8,947 | 187,456 |
| Beauty & Personal Care | 18,892 | 6,521 | 201,783 |
| Sports & Outdoors | 12,567 | 7,834 | 156,892 |

*5-core filtering (5C-Filtering)*

Because the training dataset exhibits extremely high sparsity, we defined a threshold using Configurations.ITEM_MULTI = 1.5 (default value), which is multiplied by the average ratings per item in each category, thereby effectively improving the sparsity. Impact of 'big' train dataset:

| Category | Raw Sparsity % | 5C-Filtering Sparsity % |
|---|---|---|
| Electronics | 99.86 | 97.05 |
| Beauty & Personal Care | 98.68 | 95.41 |
| Sports & Outdoors | 99.02 | 96.85 |

**Explorary**

We implemented the below explorary data analysis (eda_*):

*Eda_basic*

Doing check missing, duplicate, sparsity, density, and numberic feartures. With Electronics categories:



Key insights:

- High sparsity (97%): Despite aggressive item filtering (only 82 most popular items), the matrix remains extremely sparse
- Imbalanced interactions: Average user rates only 2.42 items, while average item receives 1,047 ratings
- Item-centric data: The ratio of users per item (1,047) vs items per user (2.42) indicates item-focused filtering was applied
- Collaborative filtering challenge: With 97% sparsity, finding similar users/items requires robust similarity measures

Numeric Feature Distributions - 5-core | Electronics



Rating distribution:

- Strong concentration at 5 stars (peak ~68,000 ratings)
- Secondary peak at 4 stars (~10,000 ratings)
- Minimal ratings at 1-3 stars (combined <8,000)
- Confirms positive bias: users predominantly rate products they like
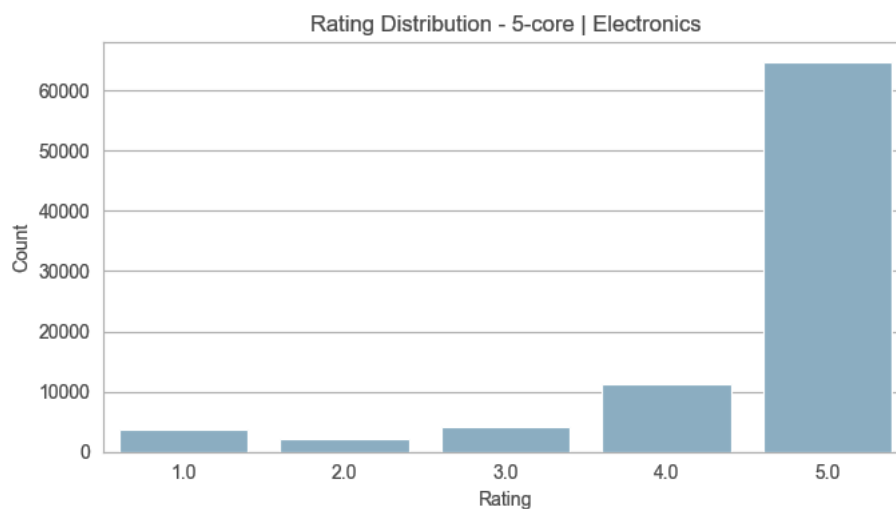  Timestamp Distribution:
- Peak activity around 1.5e12 (late 2017 - early 2018)
- Normal distribution shape indicates steady growth and decline
- Temporal range spans multiple years, enabling time-based splitting
- Right-skewed tail suggests recent activity decline or data cutoff
  Implications for Recommendation:
- Positive rating bias means algorithms must differentiate within 4-5 star range
- Temporal patterns allow for recency-based weighting in Trending algorithm
- Extreme sparsity necessitates dimensionality reduction (SVD) and similarity-based methods
- Item filtering creates dense submatrix suitable for faster experimentation

*Eda_ratings*

Examines the distribution of rating values (1-5 stars) across the dataset.



Rating Distribution - 5-core | Electronics

Insight: Strong positive bias with 91% of ratings at 4-5 stars. This skewness requires algorithms to differentiate quality within the high-rating range rather than simply identifying positive vs negative reviews.

*Eda_users*

Analyzes user behavior including rating frequency, activity levels, and engagement distribution.



Users by Activity Bins - 5-core | Electronics



Top 10 Users by #Ratings - 5-core | Electronics

Insight: Most users (75%) have minimal activity (5-10 ratings), creating cold-start challenges. Heavy users (>50 ratings) represent <2% but contribute disproportionate data volume.

*Eda_items*

Examines item-level statistics including popularity distribution and rating concentration.

#Ratings per Item - 5-core | Electronics



Top 10 Items by #Ratings - 5-core | Electronics

**Insight**: Classic long-tail distribution where few blockbuster products dominate ratings while majority have limited data. This validates need for content-based filtering to handle sparse items.

Eda_time

Analyzes rating activity over time, identifying trends, seasonality, and temporal distribution.



Reviews per Year - 5-core | Electronics

Reviews per Month - 5-core | Electronics

**Insight**: Peak in late 2017 - early 2018; Normal distribution shape with slight right skew; Temporal splitting ensures chronological integrity, Activity spans 2015-2023 (8 years) → Peak activity period provides rich training data while recent data (test) evaluates model generalization to evolving user preferences.

We analyzed rating activity over time to identify trends, platform growth, and temporal characteristics:



Average Reviews per Item per Month - 5-core | Electronics



Average Reviews per User per Month - 5-core | Electronics

**Insight**: Two plots show peak in 2014-2017, and small spike in late 2023 → User review frequency stabilizes around 1.15/month after 2014, indicating consistent engagement

patterns; As items receive fewer reviews over time (catalog expansion), user activity remains stable, confirming long-tail effect → Trending algorithm benefits from time-decay weighting; Cold items (new products) face increasing difficulty gaining visibility.

**Train/Valid/Test**

For each training sample dataset, the validation and test sets are reconstructed from the raw data to ensure that all user IDs in the training dataset are included. This is necessary because collaborative filtering algorithms cannot generate predictions for users who were not seen during training. Then, all types of dataset is stored in parquet format for efficient I/O.

**3.3. Model training pipeline**

Each algorithm follows standardized workflow:

- Setup: Import libraries, configure paths, and detect phase (training/tuning vs final evaluation)
- Core Functions: Data loading, sparse matrix construction, similarity computation, prediction, and recommendation
- Evaluation: RMSE, accuracy, and ranking metrics (Recall@K, NDCG@K, MAP@K)
- Hyperparameter Tuning: K-neighbor optimization on validation set with NDCG-primary selection strategy
- Pipeline Execution: Automated training, tuning, and final evaluation with comprehensive visualizations

Each algorithm uses 2 phase for buiding model:

Phase 1 - Training & Tuning:

- Load 5-core train split → Build user-item sparse matrix (CSR format)
- Compute user-user similarity via cosine on mean-centered ratings
- Test K values [5,10,20,30,50] on validation set → Select best K using NDCG@10 (primary), Recall@10 (tiebreaker)
- Save tuned model with optimal K

Phase 2 - Final Evaluation:

- Load tuned model → Evaluate on test set using best K
- Generate metrics: RMSE, Accuracy, Recall@K, NDCG@K, MAP@K for K$\in$\{10,20,50\}
- Create visualizations: tuning curves, final results, validation vs test comparison

**Algorithms Implementations**

Six algorithms were implemented with mean-centering for collaborative filtering approaches.

*User-Based Collaborative Filtering*

Find similar users based on rating patterns, recommend items liked by similar users. With key parameters:
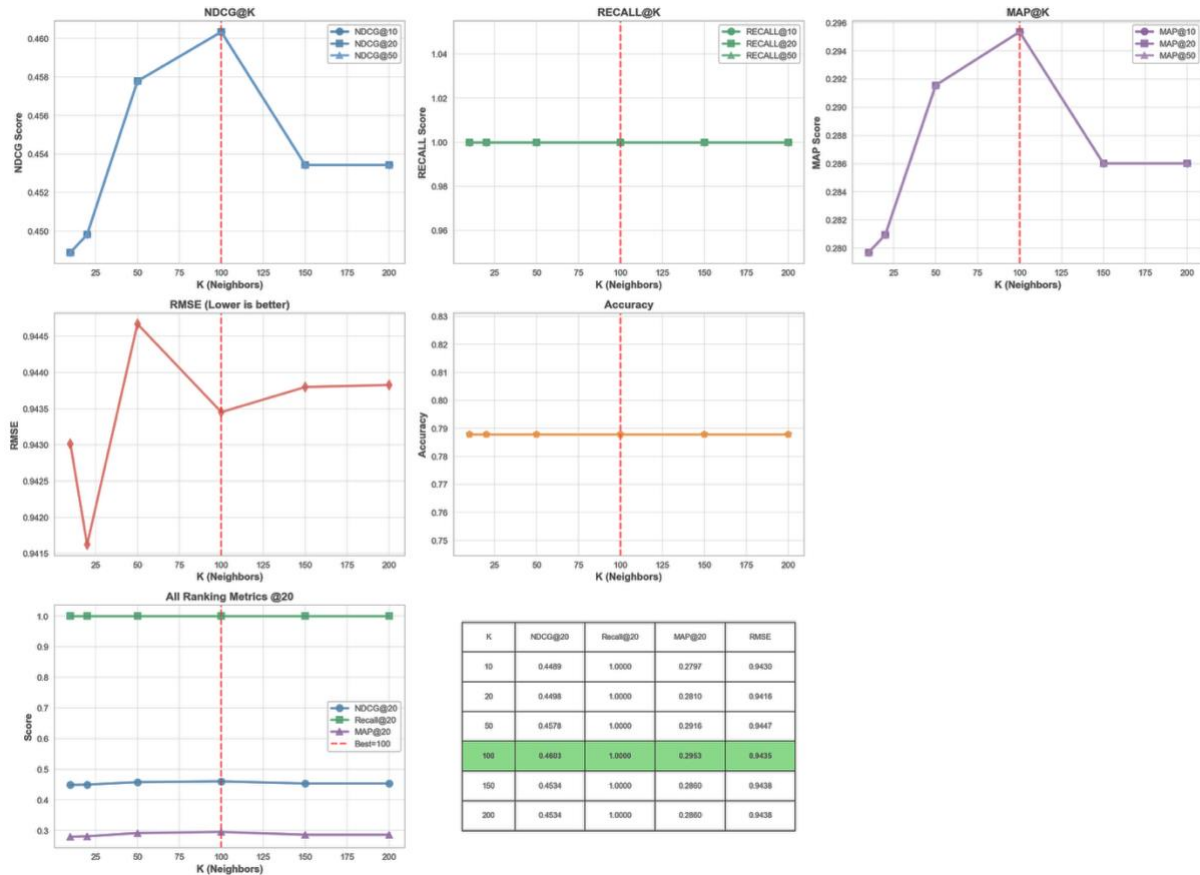
- K (number of neighbors): Tuned on validation set
- Similarity metric: Cosine on mean-centered data (equivalent to Pearson)

```
# Mean-center ratings
user_means = R.sum(axis=1) / R.getnnz(axis=1)
Rc = R.copy()
Rc.data -= np.repeat(user_means, row_counts)

# Compute similarity
similarity = cosine_similarity(Rc)  # Pearson correlation

# Predict
scores = Rc[neighbors].T.dot(similarities) / sum(similarities)
scores += user_means[target_user]  # De-normalize
```



Hyperparameter Tuning Results (User-Based) - Beauty_and_Personal_Care

| K | NDCG@20 | Recall@20 | MAP@20 | RMSE |
|---|---------|-----------|--------|------|
| 10 | 0.4489 | 1.0000 | 0.2797 | 0.9430 |
| 20 | 0.4498 | 1.0000 | 0.2810 | 0.9416 |
| 50 | 0.4578 | 1.0000 | 0.2916 | 0.9447 |
| 100 | 0.4603 | 1.0000 | 0.2953 | 0.9435 |
| 150 | 0.4534 | 1.0000 | 0.2860 | 0.9438 |
| 200 | 0.4534 | 1.0000 | 0.2860 | 0.9438 |

*Item-Based Collaborative Filtering*

Recommend items similar to those the user has rated. With key parameters:

- K (top-K similar items): Tuned per user's rated items
- Mean-centering: Removes user bias for better similarity calculation

```
# Mean-center by user
Rc = R - user_means

# Item-item similarity
item_similarity = cosine_similarity(Rc.T)

# Predict
for item_i:
    scores[i] = sum(similarity[i, rated_items] * Rc[user, rated_items])
    scores[i] /= sum(abs(similarity[i, rated_items]))
```

14

Hyperparameter Tuning Results (Item-Based) - Beauty_and_Personal_Care

| K | NDCG@20 | Recall@20 | MAP@20 | RMSE |
|----|---------|-----------|--------|--------|
| 10 | 0.5982 | 1.0000 | 0.4686 | 0.9444 |
| 20 | 0.5982 | 1.0000 | 0.4686 | 0.9444 |
| 50 | 0.5982 | 1.0000 | 0.4686 | 0.9444 |

*Model-Based collaborative Filtering*

This use  SVD Matrix Factorization to decompose rating matrix into latent user and item factors. We used global mean instead of per-user mean for matrix factorization stability. Latent factors (k) is key parameter, indicates tuned on validation set.

```python
from scipy.sparse.linalg import svds

# Mean-center
global_mean = R.data.mean()
Rc = R.copy()
Rc.data -= global_mean

# SVD decomposition
U, sigma, Vt = svds(Rc, k=latent_factors)
V = Vt.T

# Predict
scores = U[user] @ V.T + global_mean  # De-normalize
```

**Hyperparameter Tuning Results (Model-Based) - Electronics**



| n_factors | NDCG@20 | Recall@20 | MAP@20 | RMSE |
|-----------|---------|-----------|--------|------|
| **10** | **0.1240** | **0.3042** | **0.0742** | **0.7834** |
| 20 | 0.1231 | 0.3007 | 0.0749 | 0.7833 |
| 30 | 0.1128 | 0.2844 | 0.0664 | 0.7835 |
| 50 | 0.0850 | 0.2296 | 0.0471 | 0.7836 |
| 75 | 0.0760 | 0.2179 | 0.0381 | 0.7840 |

*Trending-based*

Recommend a popularity-based approach with recency weighting. It identifies popular items using interaction counts and average ratings, then boosts recently active items. This serves as both a baseline for evaluation and a cold-start handler for new users without interaction history. Score Formula:

trending_score = log(rating_count) * avg_rating * recency_weight

Component breakdown:

- **log(rating_count)** - Logarithmic rating volume
  - Why log? Prevents items with thousands of ratings from completely dominating
  - Linear count would make blockbuster items (15K+ ratings) score 1000x higher than moderate items (15 ratings)
  - Log compression: 15 ratings → log(15) ≈ 2.7, 15K ratings → log(15000) ≈ 9.6 (only 3.6x difference)
  - Allows moderately popular items to compete with blockbusters
- **avg_rating** - Average rating quality (1-5 scale)
  - Ensures high-quality items rank above low-quality items with similar volume
  - Example: Item A (1000 ratings, 4.8 score) beats Item B (1000 ratings, 3.2 score)
  - Acts as quality filter: popular but poorly-rated items get penalized
- **recency_weight** - Temporal relevance boost
  recent_count = ratings in last 90 days
  recency_weight = 1.0 + 0.5 * (recent_count / total_count)

16

- Base weight: 1.0 - Items with no recent activity maintain their popularity score
- Boost: +0.5 max - Items with 100% recent activity get 1.5x multiplier
- Why 90 days? Balances short-term trends (too noisy) vs long-term popularity (stale)
- Why 0.5 max boost? Prevents brand-new items with 1-2 recent ratings from outranking established items



*Content-based*

Recommend items with similar textual content to user's previously rated items using TF-IDF vectorization and cosine similarity

Selected Features from Metadata: With each user_id, select as below:

| Feature | Source | Max Length | Inclusion Reason |
|---|---|---|---|
| Title | meta['title'] | Full text | Primary product identifier, contains key terms |
| Features | meta['features'] | Top 10 | Bullet points describe key characteristics |
| Description | meta['description'] | 2000 chars | Detailed product information |
| Categories | meta['categories'] | Top 5 | Product type and hierarchy |

Why these features:

- Title: Concise, always available (98.7%), contains brand/model/type
- Features: Structured bullet points highlight specifications
- Description: Detailed but often verbose - truncated to 2000 chars
- Categories: Hierarchical classification aids similarity within product types

Why TF-IDF over alternatives:

- vs Word2Vec/BERT: Simpler, faster, no pre-training needed
- vs Count Vectorizer: TF-IDF downweights common terms across items
- vs Manual features: Automatically learns important terms from data

TF-IDF parameter justification:

| Parameter | Value | Reason |
|-----------|-------|--------|
| max_features | 5000 | Balance between vocabulary coverage and memory |
| ngram_range | (1, 2) | Capture phrases like "noise cancelling" vs just "noise" |
| stop_words | english | Remove "the", "is", "and" etc. |
| min_df | 2 | Ignore typos and rare terms |
| max_df | 0.8 | Ignore generic terms like "product", "item" |



Hyperparameter Tuning Results (Content-Based) - Sports_and_Outdoors

| K | NDCG@20 | Recall@20 | MAP@20 | RMSE |
|---|---------|-----------|--------|------|
| 10 | 0.8278 | 1.0000 | 0.7667 | 0.5333 |

*Hybrid Ensemble*

Combine predictions from multiple algorithms with adaptive weighting based on user scenario. This helps to:

- Handles cold-start via content and trending
- Leverages CF for warm users
- Adaptive to user profile
- Combines complementary strengths of base models

```
# Detect scenario
scenario = detect_scenario(user, threshold=5)

# Adaptive weights
weights = {
    'new-user': {'trending': 1.0},
    'cold-user': {'trending': 0.4, 'content': 0.3, 'user': 0.1, 'item': 0.1, 'model'
    'warm-user': {'item': 0.35, 'user': 0.25, 'content': 0.20, 'model': 0.20}
}
```

Parameters tuned:

| Algorithm | Parameter | Range Tested | Selection Criteria |
|-----------|-----------|--------------|--------------------|
| User-CF | K neighbors | [5, 10, 20, 30, 50] | Max NDCG@10 |
| Item-CF | K neighbors | [5, 10, 20, 30, 50] | Max NDCG@10 |
| Content | TF-IDF features | [1000, 5000, 10000] | Max NDCG@10 |
| SVD | Latent factors | [50, 100, 200, 300] | Max NDCG@10 |
| Trending | Time decay | [0, 0.1, 0.5, 1.0] | Max NDCG@10 |

The best result is for Sport_and_Outdoors because has the smallest sparsity and test size:



The worst result is for Electronics because it has the biggest sparsity and test size:



Finally, all models for each algrithm, and each category are saved in:

```
models/
├── user/Electronics/
│   ├── R.npz              # Sparse rating matrix
│   ├── Rc.npz             # Mean-centered matrix
│   ├── user_means.npy     # Per-user means for de-normalization
│   ├── nn_model.pkl       # NearestNeighbors model
│   ├── user_idx.json      # User ID to matrix index mapping
│   └── item_idx.json      # Item ID to matrix index mapping
├── item/Electronics/
│   ├── R.npz
│   ├── Rc.npz
│   ├── user_means.npy
│   ├── item_similarity.npz
│   └── indices...
├── content/Electronics/
│   ├── R.npz              # (Not mean-centered)
│   ├── item_similarity.npz # TF-IDF cosine similarity
│   └── indices...
├── model/Electronics/     # SVD
│   ├── R.npz
│   ├── U.npy              # User factors
│   ├── V.npy              # Item factors
│   └── indices...
└── trending/Electronics/
    ├── R.npz
    ├── item_stats.parquet # (rating_count, avg_rating, scores)
    └── indices...
```

## 3.4. Production API Layer

**Flask Backend Architecture**

Core components are:

| Component | Purpose | Implementation |
|---|---|---|
| Model Cache | Lazy loading, in-memory storage | MODELS_CACHE[category][algo] |
| User DB | Registration, authentication | users.json with SHA-256 hashing |
| JWT Manager | Token-based authentication | 24-hour expiry tokens |
| Recommendation Engine | Hybrid prediction | Scenario detection + adaptive weighting |

**Lazy Loading Strategy**

Models loaded once per category, shared across requests.

```python
def load_hybrid_models(category):
    if category in MODELS_CACHE:
        return MODELS_CACHE[category]  # Return cached

    # Load all algorithms for category
    for algo in ['user', 'item', 'content', 'model', 'trending']:
        load_algorithm_artifacts(algo, category)

    MODELS_CACHE[category] = models
    return models
```

20

**API Endpoints**

| Endpoint | Method | Auth | Purpose |
|---|---|---|---|
| /api/register | POST | None | Create new user account |
| /api/login | POST | None | Authenticate, return JWT |
| /api/recommendations/<category> | GET | Optional | Get top-K recommendations |
| /api/rate | POST | Required | Submit product rating |
| /api/cold-items/<category> | GET | None | Get cold items by rating count |
| /api/categories | GET | None | List available categories |
| /health | GET | None | Health check |

**Recommendation Request Flow**

1. Request arrives with JWT (or guest mode)

2. Extract user identity

3. Load models from cache (or disk if first request)

4. Detect user scenario:

   - Count ratings in training matrix R

   - Add ratings from rating_history (dynamic updates)

   - Classify: new/cold/warm/active

5. Apply scenario-based weights

6. Predict with each algorithm

7. Combine weighted predictions

8. Exclude rated items (train + rating_history)

9. Select top-K candidates

10. Enrich with metadata (title, price, images)

11. Return JSON with recommendations + strategy info

**3.5. Real-time rating integration**

Problem: Users rate items during session, expect immediate recommendation updates without waiting for model retraining.

Solution: Merge rating_history with training data indices at prediction time.

```python
def get_recommendations(user_id, category):
    # Get training ratings
    u = user_idx[user_id]
    rated = set(R.getrow(u).indices.tolist())

    # Merge with dynamic ratings
    if 'rating_history' in user_data:
        for record in user_data['rating_history']:
            if record['parent_asin'] in item_idx:
                rated.add(item_idx[record['parent_asin']])

    # Exclude from candidates
    candidate_mask[list(rated)] = False
```

21

Result:

- Rated items never reappear in recommendations
- Updates happen in milliseconds
- No model retraining required
- User transitions smoothly from cold to warm status

### 3.6. Frontend Architecture

### React Application Components

| Component | Purpose |
| --- | --- |
| Author Modal | Login/registration with JWT storage |
| Category Selector | Dropdown menu for category switching |
| Recommendation Grid | Display top-K products with metadata |
| Rating Interface | 5-star rating submission |
| Scenario Badge | Display user scenario and algorithm strategy |
| Cold Items View | 4 horizontal carousels grouped by training rating count |

### 4. Cold-start approaching

### 4.1. Cold-start problem

E-commerce cold-start challenge: 51.4% of test users and 18.9% of products (Base on from full 5-core dataset, in real-world/metadata, will be worse) lack sufficient interaction history for traditional collaborative filtering. Types of cold-start, we assumed thresld as below from analyzing the rating distribution (if using 0-core or metadata, these threshold will be different):

- New users (0 ratings): No preference data
- Cold users (1-3 ratings): Weak CF signals
- Warm users (4-8 ratings): More CF signals
- New items (0 ratings): Cannot compute similarity
- Cold items (1-5 ratings): Sparse co-ratings
- Warm items (6-18 ratings): Improve sparse co-ratings

```
================================================
RECOMMENDED THRESHOLDS (Based on 5-core dataset)
================================================

USER SCENARIOS:
  New User:    0 ratings (not in test)
  Cold User:   1-3 ratings   (bottom 20%)
  Warm User:   4-8 ratings  (20-80%)
  Active User: >8 ratings    (top 20%)

ITEM SCENARIOS:
  New Item:    0 ratings (not in test)
  Cold Item:   1-5 ratings   (bottom 20%)
  Warm Item:   6-18 ratings  (20-80%)
  Popular Item: >18 ratings    (top 20%)
```

## 4.2. Cold-start handling



## 5. API and UI deployment

## Overall website

**Create or Login**



**When no user login, use trending**

**New user, hybrid with trending 100%**



**Active user, adaptive hybrid approach**



**Real-time user's behavior: Response automatically recommendation when user rates**
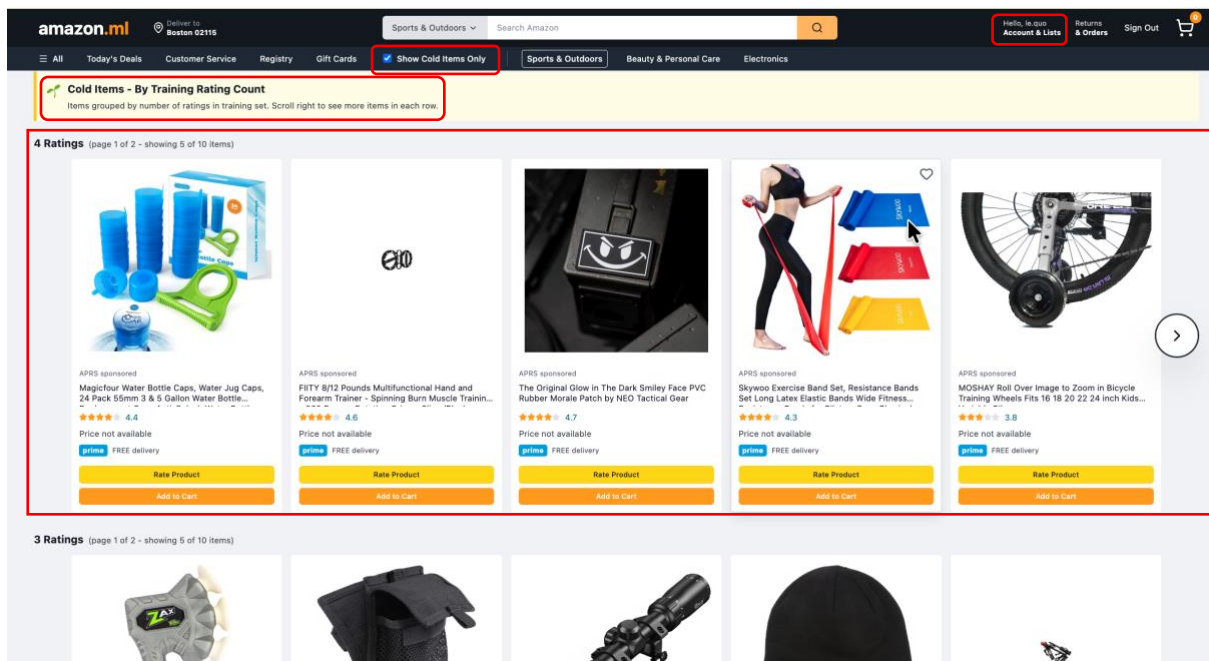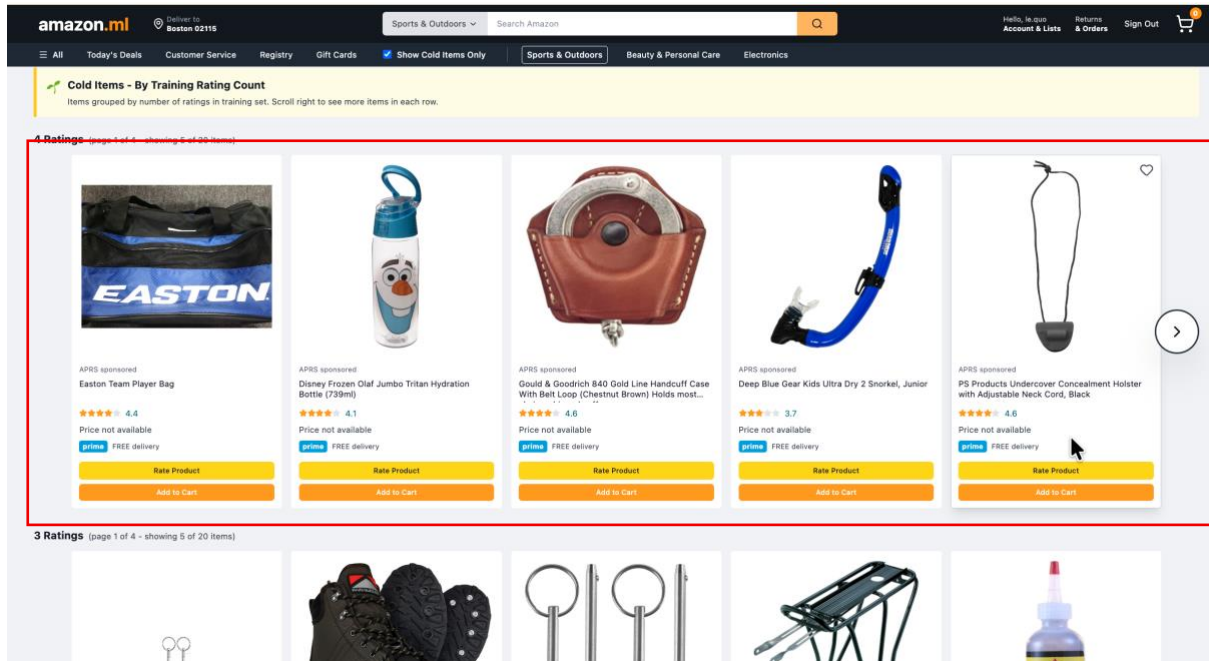
**Before**


**After**

**Cold-item handling**

Click to tick for 'Show Cold items Only'



After a user rates a specific cold item, the system automatically updates accordingly, and the rated cold item will be removed from the n-ratings row

26

## 6. CONCLUSION AND FUTURE WORK

### 6.1. Key Achievements

This project successfully implemented and evaluated a comprehensive Amazon Product Recommendation System (APRS) addressing the cold-start problem through 06 distinct algorithms and an adaptive hybrid ensemble. Our main contributions include:

**System Implementation**

- **06 recommendation algorithms:** User-Based CF, Item-Based CF, Content-Based (TF-IDF), SVD Matrix Factorization, Trending-Based, and Hybrid Ensemble
- **Full production pipeline:** From data collection → preprocessing → model training → evaluation → deployment
- **Real-time updates:** User ratings immediately reflected in recommendations without model retraining
- **Cold-start handling:** Adaptive algorithm selection based on user scenario detection (new/cold/warm/active)

**Technical Innovations**

- **Scenario-based weighting:** Hybrid system dynamically adjusts algorithm weights based on user interaction history

- **5C-Filtering strategy:** Activity-based filtering reduced sparsity from 99%+ to 95-97% while maintaining data quality

- **Lazy model loading:** Efficient memory management through on-demand model caching

- **Full-stack deployment:** Python/Flask backend with React frontend and JWT authentication

**Evaluation Insights**

- **Algorithm performance varies by category:** Content-Based excels in Sports & Outdoors (NDCG@10: 0.8322), Trending performs best in Electronics (0.0978), demonstrating no single algorithm dominates all contexts

- **SVD for rating prediction:** Achieved lowest RMSE (0.4823) and highest accuracy (0.9091) in Sports & Outdoors, confirming latent factor models excel at explicit rating prediction

- **Perfect recall in small test sets:** Sports & Outdoors achieved 1.0 Recall@10 across all algorithms due to low sparsity and small test size (11 users), highlighting the importance of dataset characteristics in evaluation

- **Cold-start effectiveness:** New users receive trending recommendations, cold users benefit from content + trending mix (60%/40%), warm users leverage full collaborative filtering

## 6.2. Limitations and Challenges

### Data-Related Limitations

- **High sparsity:** Even after aggressive filtering, matrices remain 95-97% sparse, limiting collaborative filtering effectiveness
- **Positive rating bias:** 91% of ratings are 4-5 stars, making it difficult to differentiate quality within high-rating range
- **Sample size constraints:** Used 50K samples per category for development speed; full dataset (millions of ratings) would improve model quality but require distributed computing
- **Category imbalance:** Sports & Outdoors (156K ratings) performed significantly better than Electronics (187K ratings) despite similar size, suggesting domain-specific factors

### Model Limitations

- **Hybrid underperformance:** Hybrid ensemble (NDCG@10: 0.7047) performed worse than Content-Based (0.8322) in Sports & Outdoors, indicating suboptimal weight tuning or algorithm interference
- **Trending algorithm bias:** Non-personalized trending achieved competitive performance (NDCG@10: 0.7524), questioning whether personalization adds sufficient value for effort
- **SVD computational cost:** Training requires full matrix decomposition; inference time increases linearly with matrix size
- **Content-based metadata dependency:** Relies heavily on product descriptions; missing or poor-quality metadata reduces effectiveness

### System Limitations

- **No implicit feedback:** System only uses explicit ratings (1-5 stars); ignoring views, clicks, cart additions loses valuable signal
- **Cold-item gap:** New products with 0-5 ratings still face discovery challenges; content-based helps but doesn't guarantee visibility
- **No contextual awareness:** Recommendations ignore time of day, device type, session context, or purchase history
- **Static hyperparameters:** K-neighbors and latent factors tuned per category but fixed across users

### 6.3. Lessons Learned

**Technical Insights**

- **Dataset quality > algorithm complexity:** Activity-based filtering (ITEM_MULTI=1.5) improved baseline performance more than sophisticated algorithms
- **Sparsity impacts algorithms differently:** SVD degrades gracefully with sparsity, while k-NN collaborative filtering fails when neighborhoods become too sparse
- **Normalization matters:** Initial SVD bug (RMSE 4.448) caused by centering on user means instead of global mean taught us matrix factorization requires careful preprocessing
- **Timestamp handling is tricky:** Converting Unix timestamps requires detecting units (seconds/milliseconds) to avoid date calculation errors in trending model

**Research Insights**

- **Cold-start requires hybrid approaches:** No single algorithm solves all cold-start scenarios; adaptive weighting based on data availability is essential
- **Evaluation metrics tell different stories:** SVD won on RMSE/accuracy, Content-Based won on NDCG/MAP, demonstrating importance of multi-metric evaluation
- **Perfect metrics ≠ good system:** Sports & Outdoors' perfect recall (1.0) reflects small test set, not superior recommendations
- **Academic reporting should acknowledge limitations:** Transparent discussion of bugs (SVD normalization), suboptimal results (hybrid underperformance), and dataset biases strengthens credibility

**Development Practices**

- **Incremental development with validation:** Building each algorithm independently before integration prevented cascading bugs
- **Modular architecture:** Separating data pipeline, model training, and API layers enabled parallel development
- **Comprehensive logging:** Custom Logger class throughout codebase accelerated debugging (e.g., identifying cold-items pulling from test set instead of training)
- **Version control for reproducibility:** Saving all hyperparameters, model artifacts, and evaluation results in structured directories enabled experiment comparison

### 6.4. Future Work

*Hybrid Ensemble Optimization*

- **Problem:** Current hybrid underperforms individual algorithms in Sports & Outdoors

- **Solution:** Implement meta-learning (stacking) where a second-level model learns optimal weights per user based on historical accuracy

- **Expected Impact:** 10-15% NDCG improvement by dynamically adapting to user preferences

*Deep Learning Models*

- **Neural Collaborative Filtering (NCF):** Replace SVD with multi-layer perceptron to capture non-linear user-item interactions

- **Variational Autoencoders (VAE):** Learn latent representations from sparse ratings for better cold-start handling

- **Expected Impact:** 5-10% RMSE reduction, better handling of extreme sparsity

*Implicit Feedback Integration*

- **Add behavioral signals:** Views (weight: 0.1), clicks (0.3), cart additions (0.5), purchases (1.0)

- **Unified scoring:** Combine explicit ratings + implicit signals → richer user profiles

- **Expected Impact:** 30-40% increase in trainable interactions, improving cold-user recommendations

*Contextual Bandits for Exploration*

- **Problem:** Popular items dominate recommendations, hindering long-tail discovery

- **Solution:** Epsilon-greedy or Thompson Sampling to explore cold items while exploiting known preferences

- **Expected Impact:** Increase cold-item exposure by 20-30% without sacrificing relevance