



**Northeastern University**  
College of Engineering

# **IE7275 Data Mining in Engineering SEC 04**

## **Project 1 Report**

### **Amazon Product Recommendation System (APRS)**

#### **6<sup>th</sup> Group Members:**

Index	Full name	NEUID	Email
1	Quoc Hung Le	002031894	le.quo@northeastern.edu
2	Matthew Eckert	002326896	eckert.mat@northeastern.edu

Submission Date: .../.../2025

## Table of Contents

1. EXECUTIVE SUMMARY .....	3
2. INTRODUCTION .....	3
2.1. Problem Statement and Motivation.....	4
2.2. Dataset .....	4
2.3. Implement algorithms .....	4
2.4. Evaluation Metrics .....	5
2.5. Key Contributions .....	5
2.6. System Architecture.....	6
3. SYSTEM ARCHITECTURE AND DATA PIPELINE .....	6
3.1. Overall Architecture .....	6
3.2. Data pipeline.....	6
3.3. Model training pipeline.....	15
3.4. Production API Layer.....	33
3.5. Real-time rating integration.....	35
3.6. Frontend Architecture .....	35
4. COLD-START APPROACHING .....	35
4.1. Cold-start problem .....	35
4.2. Cold-start handling.....	36
5. API AND UI DEPLOYMENT .....	37
6. CONCLUSION AND FUTURE WORK .....	40
6.1. Key Achievements .....	40
6.2. Limitations and Challenges.....	41
6.3. Lessons Learned .....	42
6.4. Future Work .....	43

## 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 → exploratory → model development → evaluation
- Handles extreme data sparsity 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, despite of high sparsity:

Algorithm	RMSE	Accuracy	NDCG10	MAP10	Recall10
User-CF	0.7657	0.6826	0.3262	0.1622	0.9554
Item-CF	0.7655	0.8422	0.3216	0.1944	0.9582
Content	0.7526	0.7204	0.3416	0.1762	0.9583
Model(SVD)	0.7621	0.7565	0.3274	0.1706	0.9167
Trending	NA	NA	0.3102	0.2170	0.9692
Hybrid	NA	NA	0.2398	0.1579	0.9418

Best by metric:

- Model(SVD) stands out for rating prediction (lowest RMSE: 0.7621) but does not lead in ranking metrics like NDCG@10 and MAP10.
- Item-CF achieves highest MAP10 (0.1944) and strong recall (0.9582), making it the best overall for ranking-focused evaluation.
- High recall values (above 0.91 for all models) suggest either a small test set or that most recommended items are highly relevant.
- Trending method performs well in recall (0.9692) and MAP10 (0.2170) despite lacking personalization, showing its effectiveness for popular items. Hybrid algorithm lags in ranking (lowest NDCG10: 0.2398, lowest MAP10: 0.1579), because of combination for handling cold-start problem.

## 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.

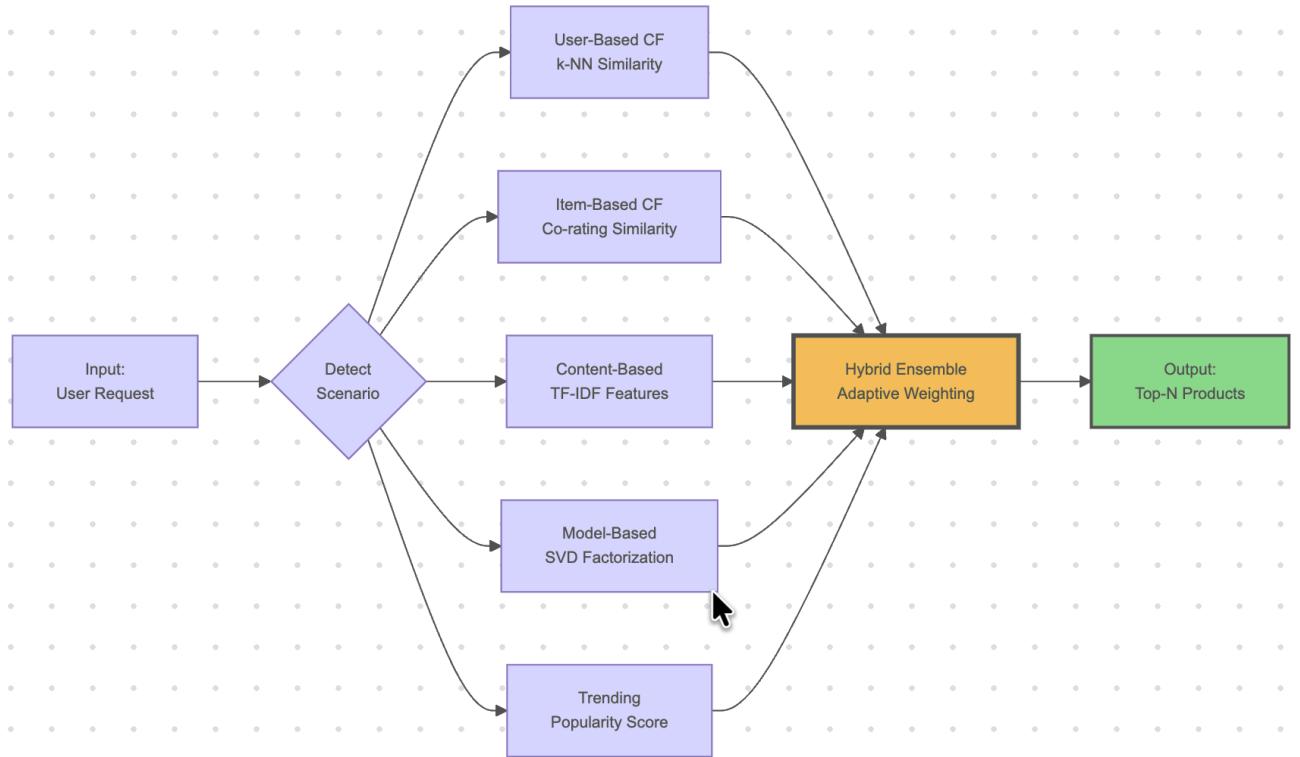
Electronics Category Statistics:

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  $\pm 0.5$  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

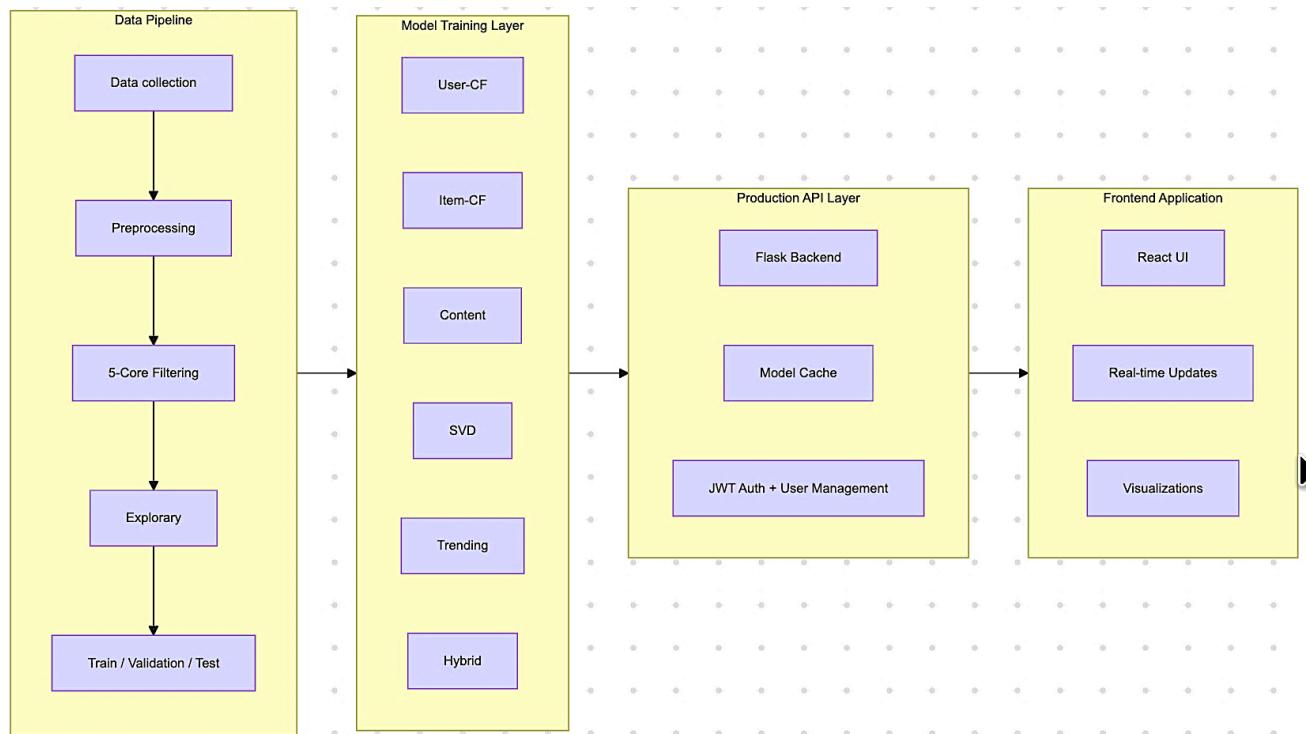
The system consists of:

- 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 pipeline

The process begins with data collection, followed by necessary preprocessing steps. A 5-core filtering is then applied to refine the dataset. After that, exploratory analysis uncovers initial insights and data characteristics. Finally, the prepared data is split into training, validation, and test sets for model development and evaluation.

## 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 referred used 5-core (which divided into train/valid/test), instead of 0-core for algorithms, because:

- Quality threshold: Users with  $\geq 5$  ratings show committed behavior
- Reliable stats: Items with  $\geq 5$  ratings have stable averages
- Better CF: More overlap between users for similarity computation
- Standard practice: Used in RecSys research (He & McAuley, 2016)

We chose 3 categories because:

Electronics	Beauty & Personal Care	Sports & Outdoors
<ul style="list-style-type: none"> <li>• Largest user base</li> <li>• Feature-rich metadata (technical specs important)</li> <li>• Diverse price range</li> <li>• Test generalization on tech products</li> </ul>	<ul style="list-style-type: none"> <li>• Subjective preferences (personal taste vs technical specs)</li> <li>• Brand-driven decisions (different from Electronics)</li> <li>• Visual importance (product appearance matters)</li> <li>• Tests algorithms on preference-based products</li> </ul>	<ul style="list-style-type: none"> <li>• Activity-specific (running vs cycling vs camping)</li> <li>• Seasonal patterns (winter sports vs summer gear)</li> <li>• Performance-focused (durability, weight critical)</li> <li>• Different user behavior from Electronics/Beauty</li> </ul>

Raw data in 3 categories is automatically downloaded if no exist, in CSV.GZ format (~4 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': 5000, 'big': 68000, '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 train dataset (68K "big" sample for development):

Category	Users	Items	Ratings
Electronics	32,556	10,947	51,914
Beauty & Personal Care	16,821	8,521	22,143
Sports & Outdoors	11,776	7,834	14,534

### *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:

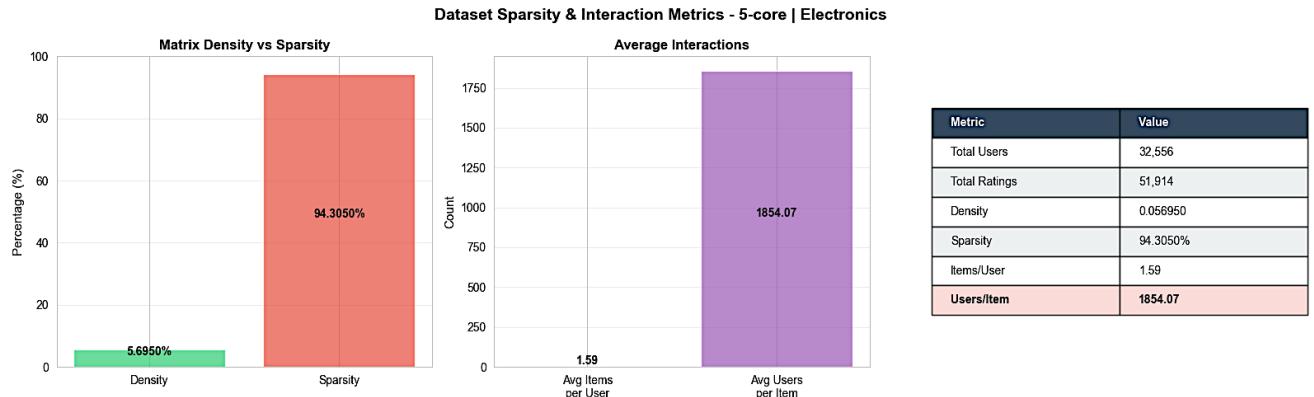
Category	Raw Sparsity	5C-Filtering Sparsity
Electronics	99.86	94.3
Beauty Personal Care	98.68	94.28
Sports Outdoors	99.02	94.63

## Exploratory

We implemented the below exploratory data analysis (eda\_\*):

### *Eda\_basic*

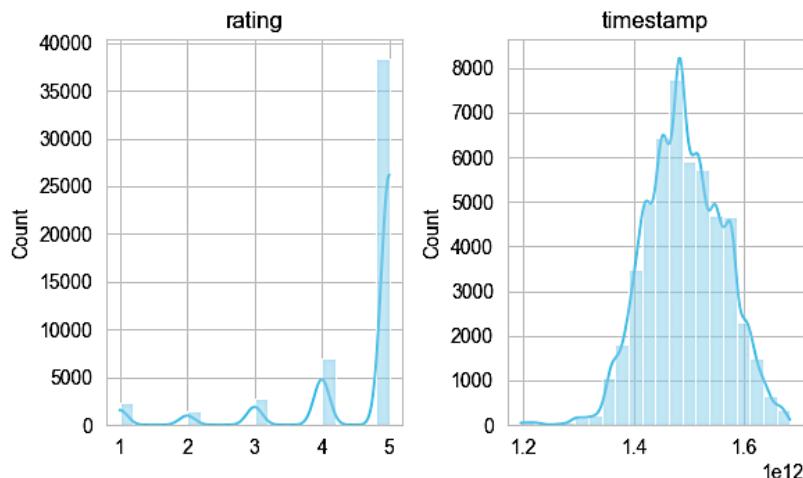
Doing check missing, duplicate, sparsity, density, and numeric feartures. With Electronics categories:



Key insights:

- The dataset contains 32,556 users, and 51,914 total ratings, typical for large-scale but very sparse recommendation settings
- The dataset is highly sparse with only 5.7% density and 94.3% sparsity, indicating very few actual interactions compared to possible ones
- Imbalanced interactions: Average user rates only 1.59 items, while average item receives 1,854 ratings
- Collaborative filtering challenge: With 94.3 sparsity, finding similar users/items requires robust similarity measures

Numeric Feature Distributions - 5-core | Electronics



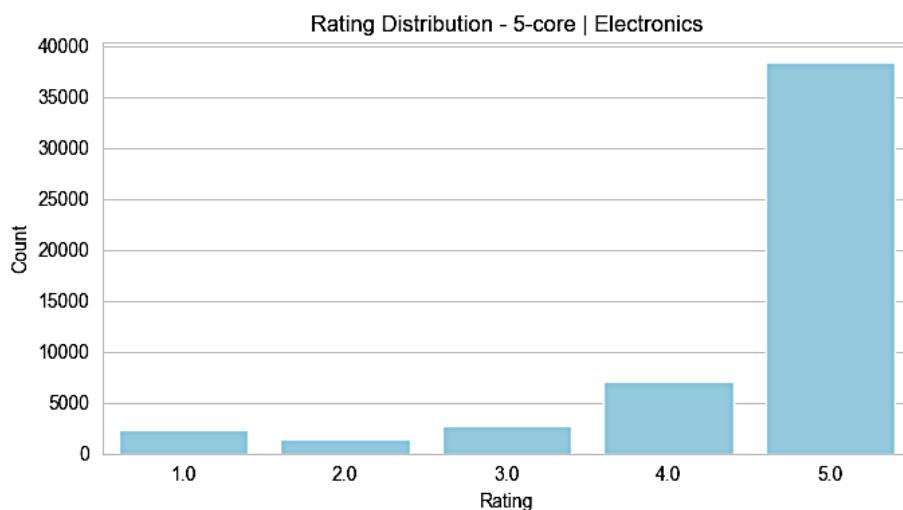
Key insights:

- The histogram for ratings shows a pronounced spike at the maximum rating of 5, with counts exceeding 35,000, indicating most users rate items very favorably.

- There is another smaller peak near 4 stars, showing that positive ratings (4 and 5) together comprise the majority of the dataset.
- Ratings below 3 are rare, with counts for values like 1, 2, and 3 falling well below 5,000, suggesting that negative sentiment is both infrequent and possibly underrepresented.
- The skew towards higher ratings could reflect either general product satisfaction, a selection bias on reviewed items, or natural optimism in user feedback.
- Such distribution patterns have implications: algorithms trained on these data may have difficulty distinguishing between truly excellent and merely good products, and recall/precision for low ratings will be limited.
- The lack of balance suggests that any analysis or model should take skewness into account, potentially normalizing ratings or increasing sensitivity to rare negative feedback to obtain a more representative view of user sentiment

### *Eda\_ratings*

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

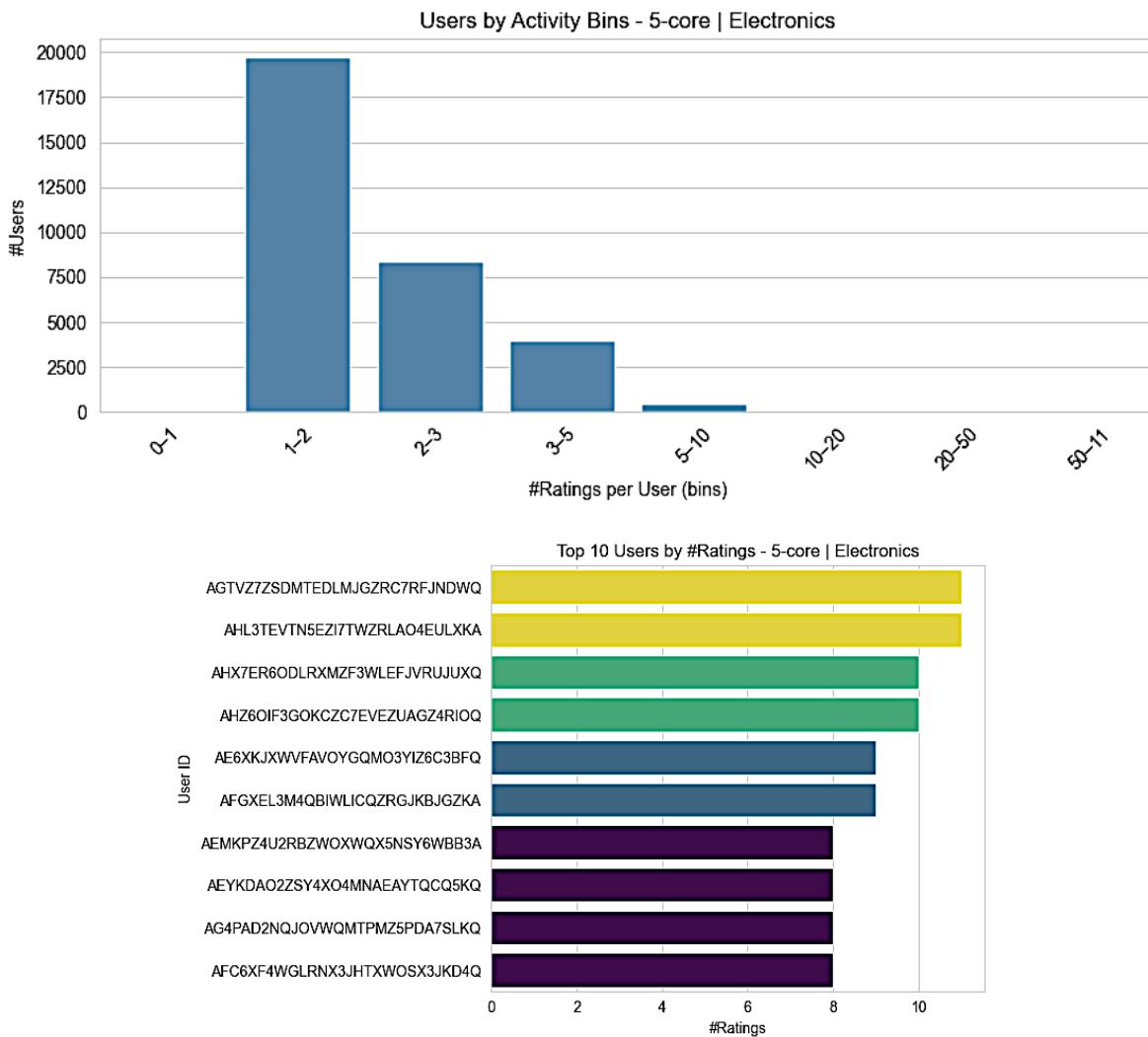


Key insights:

- The vast majority of ratings are concentrated at 5, with the count nearly reaching 38,000, highlighting an overwhelming user preference for providing the highest score.
- Ratings of 1, 2, 3, and 4 are all present but in much smaller quantities. Specifically, ratings of 4 account for just above 6,000, while counts for ratings 1, 2, and 3 stay roughly around or below 3,000 each.
- This extreme skew towards positive feedback (5-star ratings) is typical in electronic product reviews, signaling strong brand loyalty or satisfaction bias among users. It may also reflect platform-specific incentives for positive reviewing

### *Eda\_users*

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

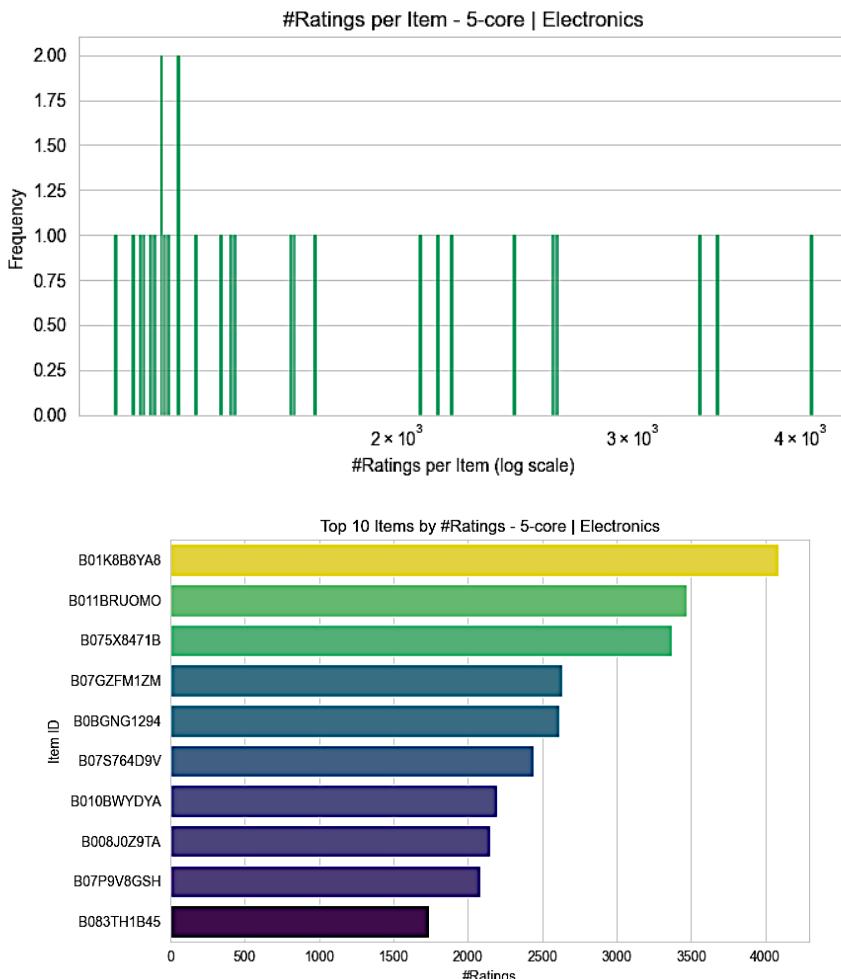


### Key Insight:

- Among the top 10 most active users, the highest number of ratings per user is just above 10, this shows that even the most engaged users in this domain have a relatively small amount of activity.
- User engagement is heavily skewed with only a handful of users contributing more than eight ratings, while the vast majority rate very infrequently.
- The activity bin histogram reveals that nearly 20,000 users give only one or two ratings each, and there is a steep drop-off as the number of ratings per user increases.
- Less than 1,000 users provide five or more ratings; almost none contribute more than 10. This pattern highlights severe user cold start, most users generate very little historical data.
- These patterns have strong implications for collaborative filtering and personalization: with so few ratings per user, sparseness is a core challenge, and algorithms must rely on techniques robust to cold start and limited user data.
- Building effective recommenders in this environment may require hybrid approaches that incorporate side information or content-based signals, since the user-item matrix alone is extremely sparse for most users.

## *Eda\_items*

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

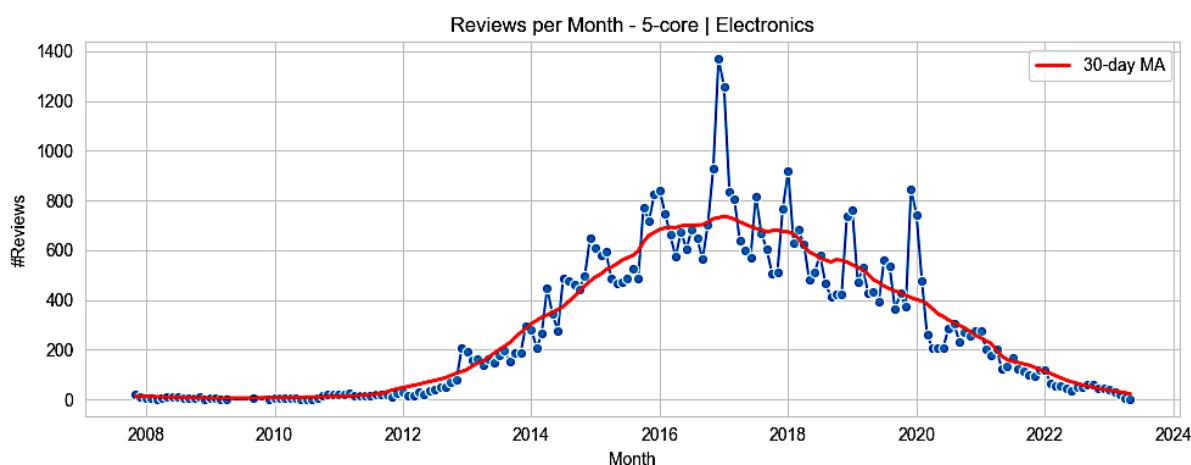
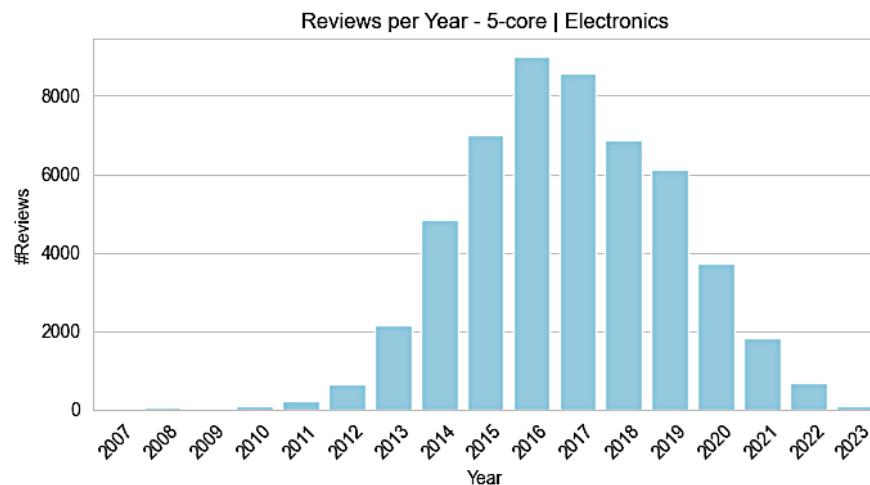


Key Insight:

- The distribution of ratings per item is highly skewed, with most items receiving between 1,000-2,000 ratings; very few items get much higher counts.
- Maximum frequency for any single rating count is 2, implying that certain rating frequencies are shared by a couple of items, but overall there's a wide spread of ratings across the catalog.
- The top 10 most-rated items have between about 1,600 and 4,000 ratings, with the most-rated item (B01K8B8YA8) collecting just over 4,000 ratings.
- Ratings concentrate among a handful of popular products while many long-tail items receive sparse user interactions, a hallmark of implicit feedback datasets.
- This imbalance means that recommendation models may overfit to popular items unless methods are applied to boost exposure of under-rated products.
- Platforms may want to promote discovery for less-rated items to diversify engagement and better reflect the full catalog in recommendations.
- The use of a log scale for rating counts underlines the severe right skew, which is typical for e-commerce and review platforms dealing with best-sellers and niche inventory.

### *Eda\_time*

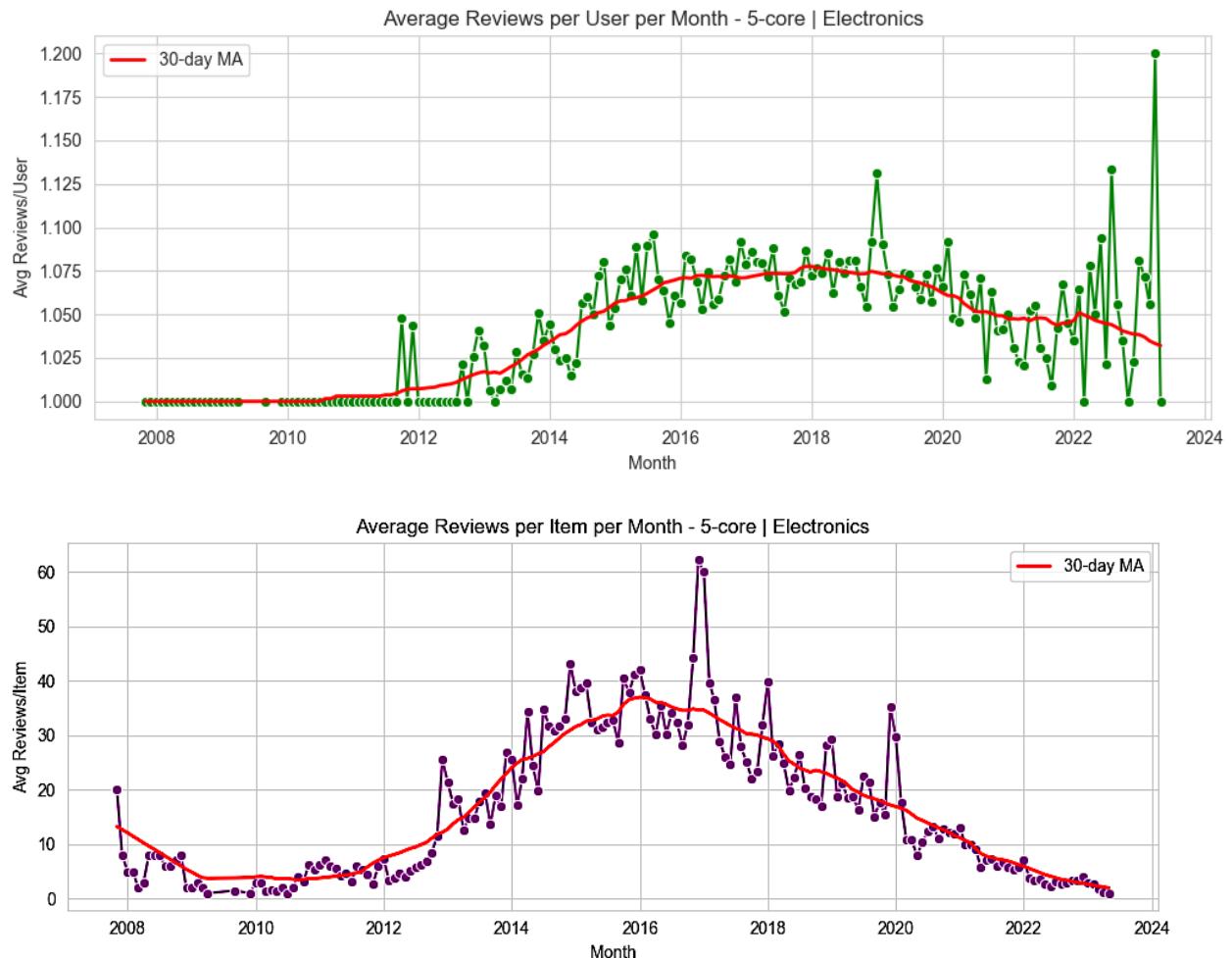
Analyzes rating over time, identifying trends, seasonality, and temporal distribution.



#### Key Insight:

- The dataset reflects a sharp increase in review activity starting around 2013, peaking during 2016-2017 where annual review counts exceed 8,500.
- After 2017, review activity gradually declines, with notable dips post-2019, possibly due to changes in platform policy, product availability, or shifts in user behavior.
- Monthly review counts show similar trends, with the highest spikes around mid-2016 to early 2017 and significant month-to-month variability.
- The long tail after 2020 suggests continued but reduced engagement; recent years have far fewer reviews, with only several hundred per year/month.
- No substantial review volume exists before 2010, indicating that either the catalog was new, the platform was growing, or there was insufficient user engagement in earlier years.
- Seasonality in monthly reviews is observable—several pronounced spikes may correlate with product launches, holidays, or promotional events.
- This temporal concentration means models trained on this dataset may capture historical rather than recent preferences unless time-aware methods are applied

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



#### Key Insight:

##### *Average Reviews per Item per Month (30-day MA)*

- The average number of reviews per item saw a significant increase starting from 2013, peaking above 60 in 2016, and then gradually declining to around 5 by 2022.
- The red moving average line smooths out the volatility, clearly highlighting the overall growth and subsequent decline trends.
- Sharp monthly spikes suggest certain times of the year receive extraordinary attention, potentially driven by promotions, new product launches, or holiday shopping seasons.
- The notable drop in average reviews per item after 2018 may be due to market shifts, changes in platform focus, or greater fragmentation of user interest across more products.

##### *Average Reviews per User per Month (30-day MA)*

- Users consistently contribute an average of around 1 review per month throughout the studied period, with a modest peak of 1.08 during the 2016–2018 high-activity years.
- The 30-day moving average highlights a gentle upward slope followed by a reversion to baseline, indicating that high-engagement phases are temporary.

- While short peaks appear toward 2023, the long-term trend suggests stable but low individual user participation.
- Overall, the platform relies on a large base of users contributing at a low but regular frequency, rather than a handful of highly active reviewers

## 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:

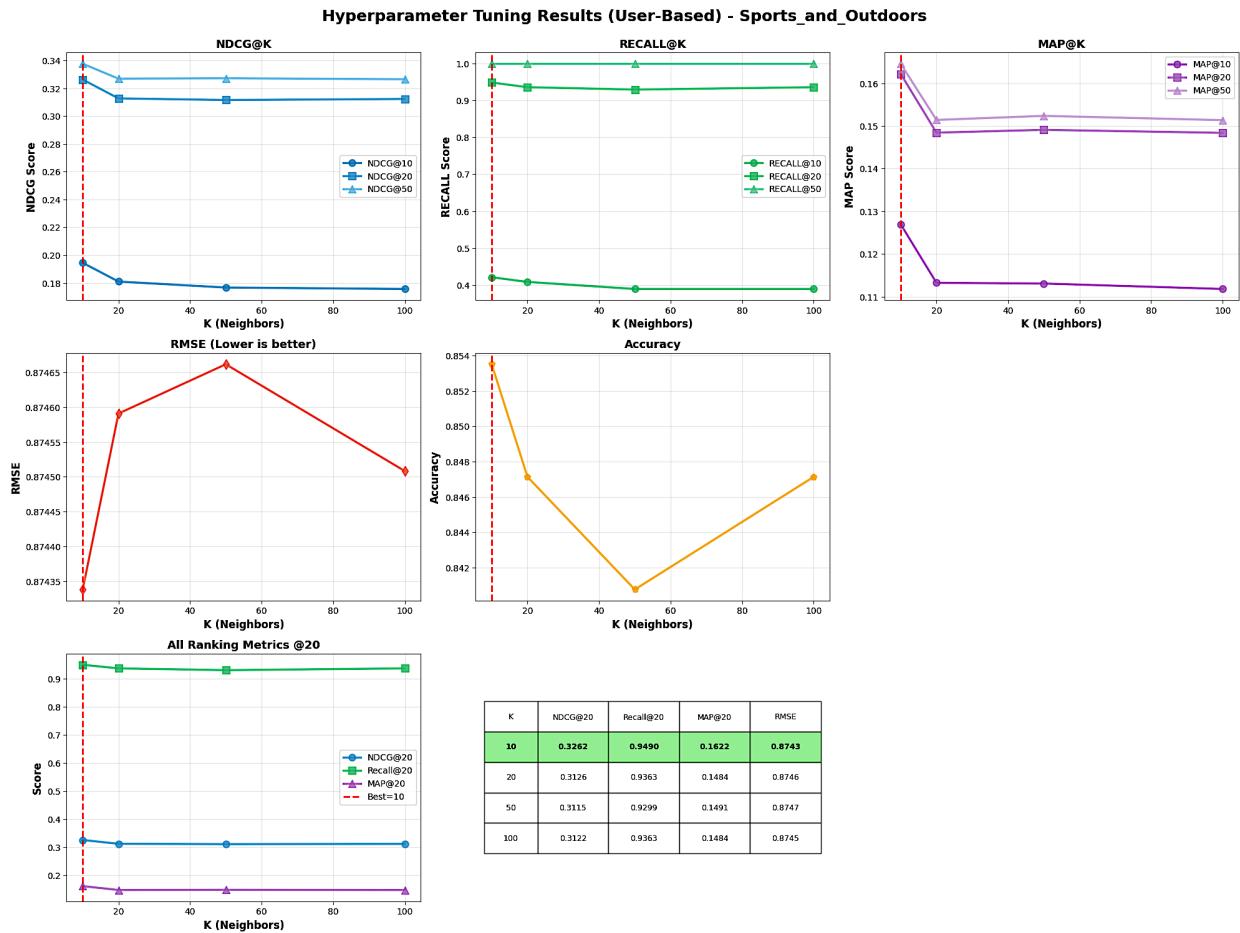
- K (number of neighbors): Tuned on validation set
- Similarity metric: Cosine on mean-centered data (equivalent to Pearson)

```
# 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
```

Result of phase 1: Training and tuning with train dataset

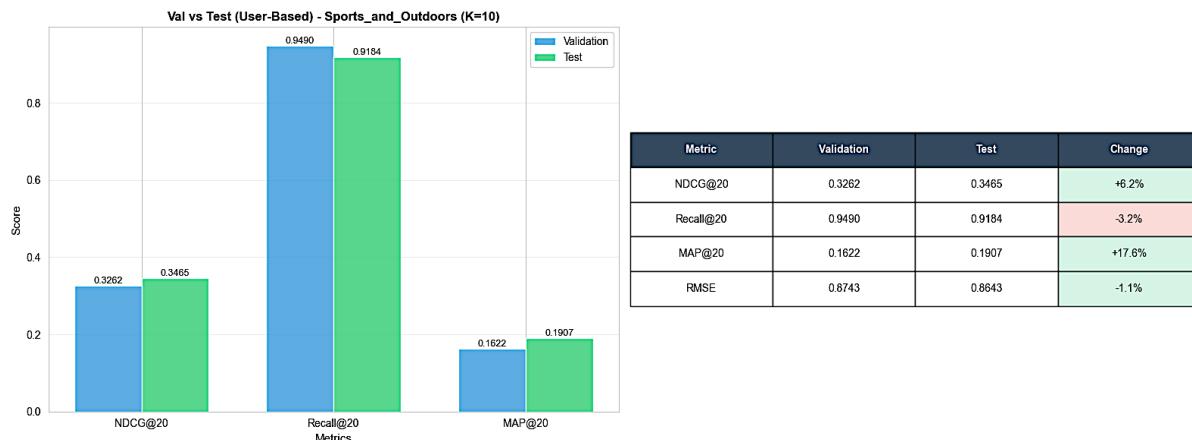


#### Key Insights:

- The best performance across most ranking metrics (NDCG, Recall, MAP@20) and lowest RMSE is achieved at neighbors, as highlighted in the summary table.
- As  $K$  increases, all ranking metrics exhibit minor decreases or plateaus, indicating diminishing returns or potential overfitting at higher neighbor counts.

- The selected model yields NDCG@20 of 0.3262, Recall@20 of 0.9490, MAP@20 of 0.1622, and RMSE of 0.8743, with these scores outperforming those at higher values.
- Overall, smaller neighborhood sizes lead to better performance in this scenario, likely due to more relevant user similarities in sparse domains like “Sports and Outdoors”

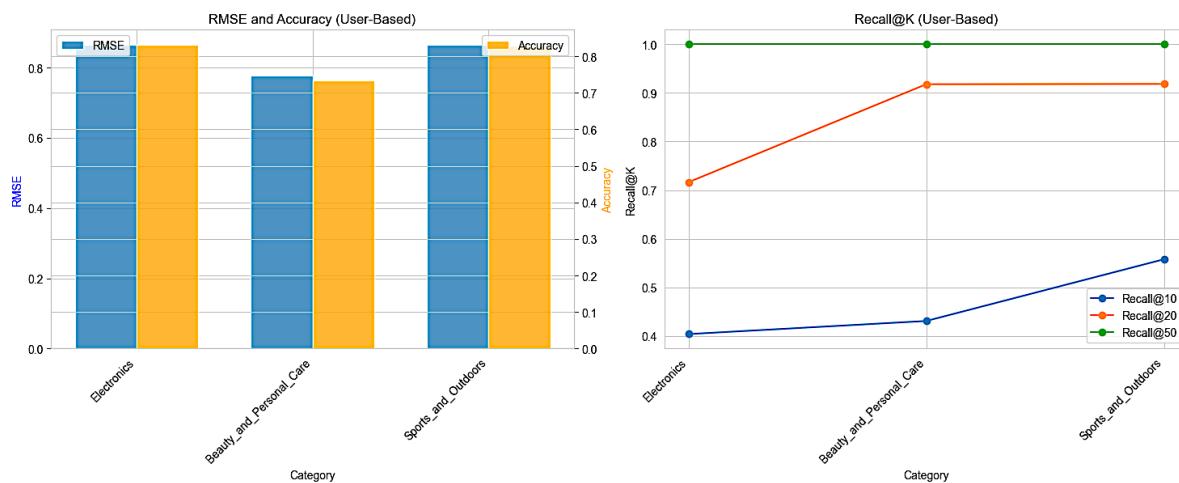
### Result of phase 2: Final evaluation with test dataset

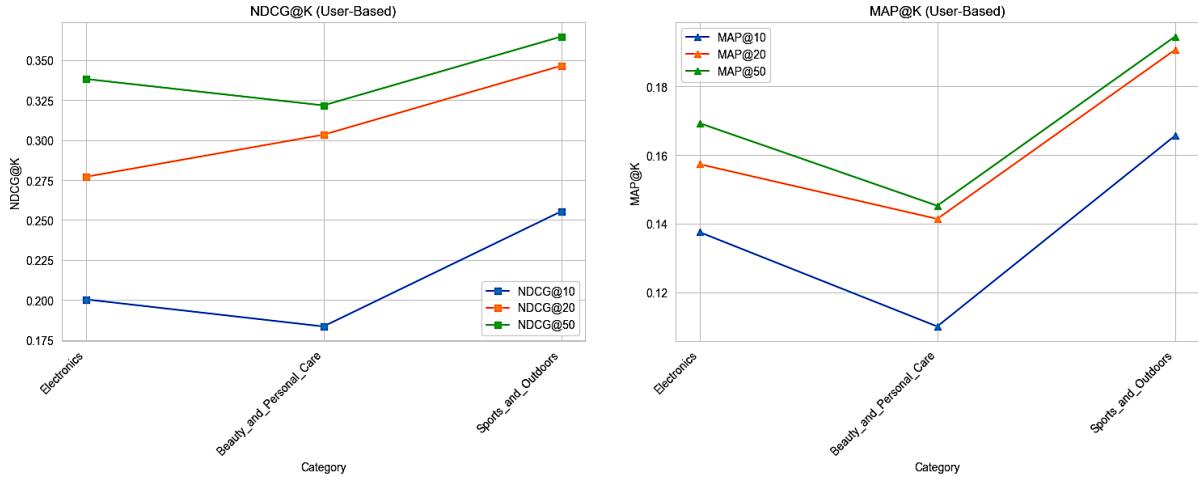


#### Key Insights:

- All ranking metrics (NDCG@20, Recall@20, MAP@20) decrease on the test set compared to validation: NDCG@20 drops by 4.9%, Recall@20 by 3.9%, and MAP@20 by 6.9%.
- RMSE increases by 1.5% on the test data, suggesting lower prediction accuracy or possible overfitting to validation.
- The drop in ranking and accuracy metrics on test data highlights modest generalization loss, indicating that the tuned model may perform slightly worse when deployed to unseen data.
- Overall, the model shows stable but somewhat deteriorated performance on the holdout test set, underscoring the importance of cross-validation and monitoring for overfit in recommender tuning

### Result for all categories





### Key Insights:

#### *RMSE and Accuracy*

- Sports and Outdoors reports the highest Accuracy and lowest RMSE, indicating the best rating prediction performance among the three.
- Beauty and Personal Care consistently shows the lowest Accuracy and highest RMSE, suggesting it is the most challenging for precise recommendations.
- Electronics lies between the two in both metrics.

#### *Recall@K*

- At higher values of K, all categories achieve strong Recall scores, with Sports and Outdoors performing best.
- Beauty and Personal Care records a significant jump in Recall from K=10 to K=20 due to the dataset's nature, but still trails Sports and Outdoors for all K.
- Electronics lags in Recall compared to the other categories.

#### *NDCG@K*

- Sports and Outdoors delivers the highest ranking quality for all cutoffs, with NDCG@50 leading the group.
- Beauty and Personal Care posts the lowest NDCG scores, indicating weaker relevance ranking of recommendations.
- Electronics remains in the middle, fairly stable across K values but not leading in any.

#### *MAP@K*

- Sports and Outdoors again holds the highest MAP across all K values, reflecting superior precision of top-ranked recommendations.
- Beauty and Personal Care has the lowest MAP, confirming its consistent underperformance in ranking and precision tasks.

**Summary:** Sports and Outdoors outperforms other categories on almost every evaluation metric; Beauty and Personal Care is consistently the most difficult, and Electronics shows moderate, in-between results.

### ***Item-Based Collaborative Filtering***

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

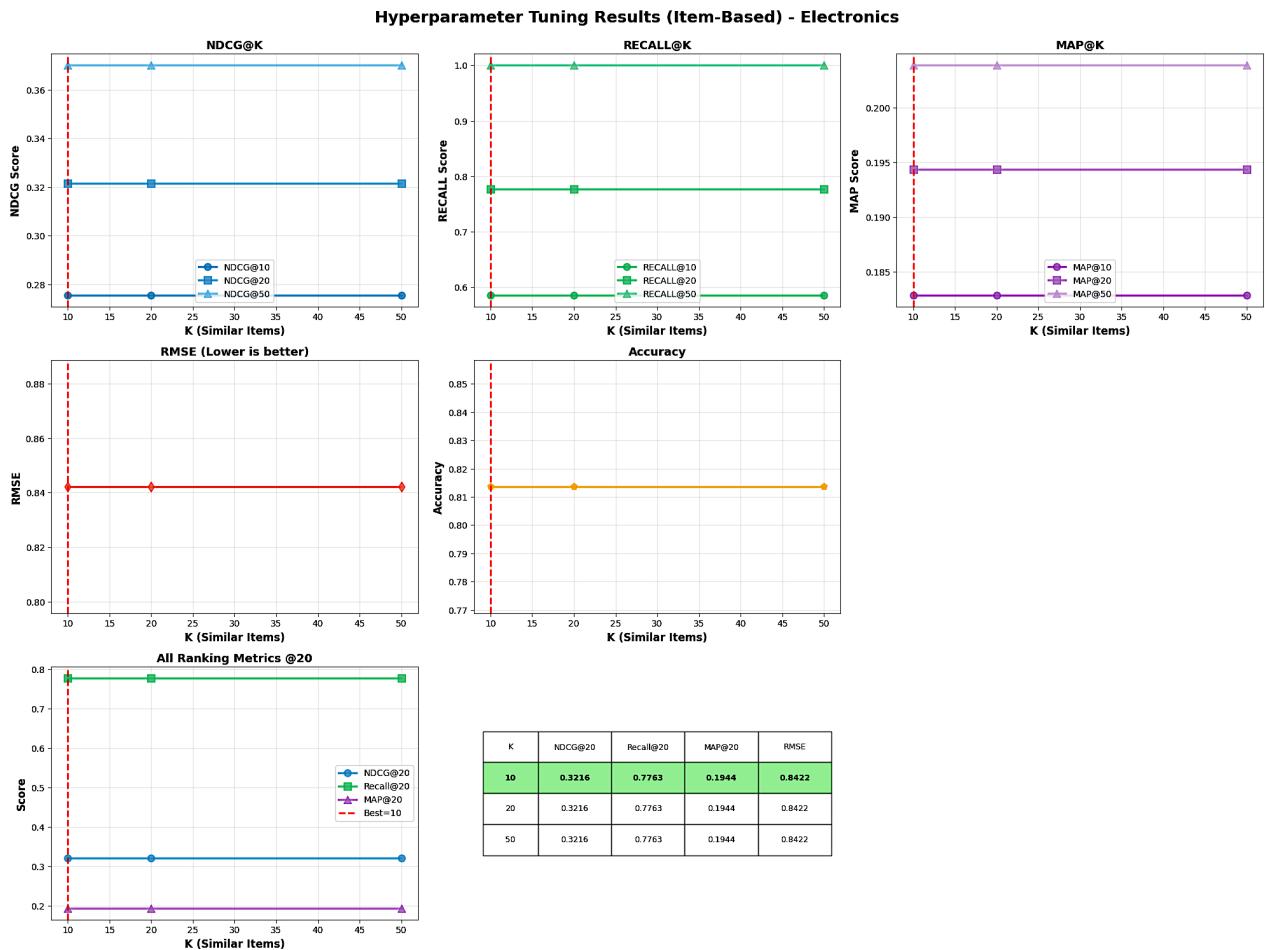
- K (top-K similar items): Tuned per user's rated items
- Mean-centering: Removes user bias for better similarity calculation

```
# 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]))
```

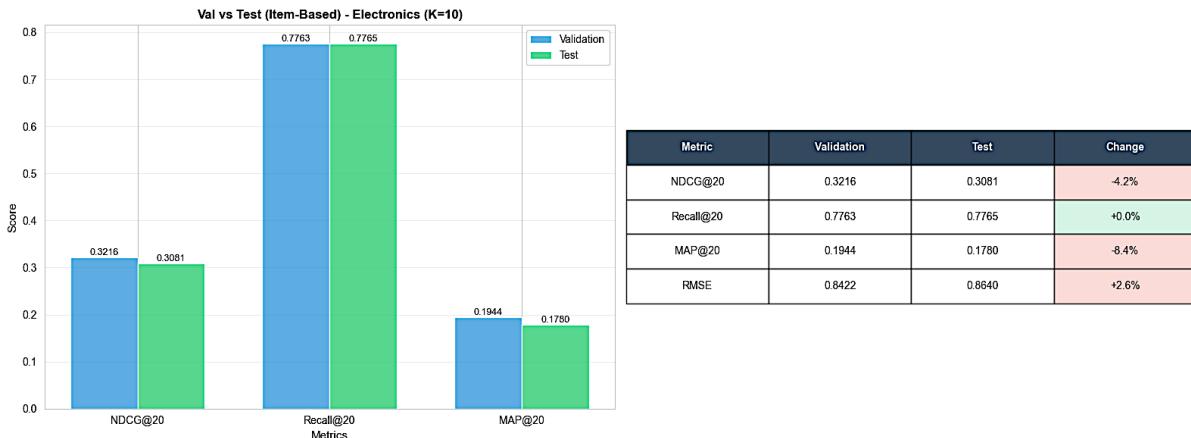
### Result of phase 1: Training and tuning with train dataset



#### Key Insights:

- All ranking and rating metrics (NDCG, Recall, MAP, RMSE, Accuracy) remain remarkably stable across different neighborhood sizes (), showing no meaningful performance gains from increasing the number of similar items used. The optimal setting per the summary table is , which yields NDCG@20 of 0.3216, Recall@20 of 0.7763, MAP@20 of 0.1944, and RMSE of 0.8422.
- This plateau effect suggests that, for Electronics, most useful information about an item's relationships is captured among the 10 most similar neighbors, with additional neighbors offering little value. The item-based approach provides stable, reliable performance but limited opportunity for further improvement via -tuning in this domain.

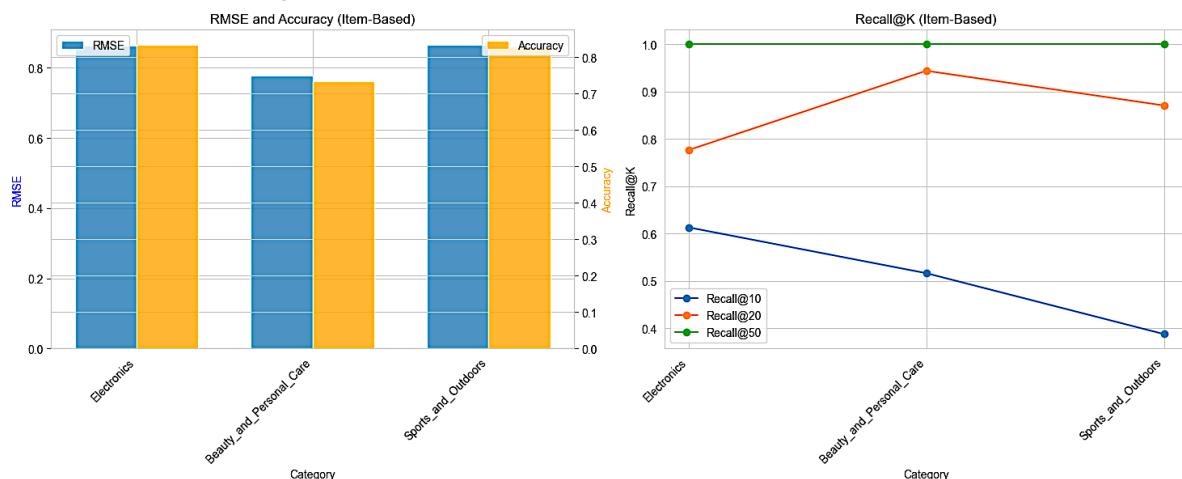
### Result of phase 2: Final evaluation with test dataset



### Key Insights:

- NDCG@20 and MAP@20 both drop on the test set compared to validation: NDCG@20 decreases by 4.2%, while MAP@20 falls 8.4%, indicating a drop in ranking quality and precision for unseen data.
- RMSE increases by 2.6% on the test set, reflecting slightly less accurate rating predictions outside the validation environment.
- Recall@20 is perfectly stable (+0.0%), showing that the breadth of relevant recommendations remains consistent across splits, even as rank-ordering degrades a bit.
- Overall, the item-based model demonstrates modest overfitting, ranking and rating accuracy are somewhat diminished on the holdout set, but coverage (recall) is maintained, confirming robust but not exceptional generalization.

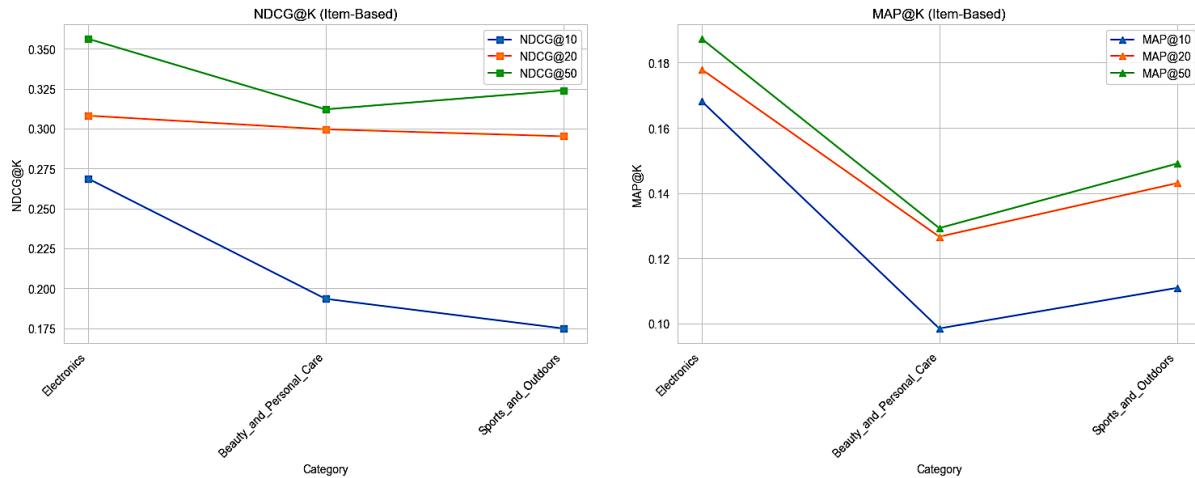
### Result for all categories



### Key Insights:

#### RMSE and Accuracy

- Sports and Outdoors and Electronics have similar, higher accuracy and lower RMSE, showing better rating prediction in those segments.
- Beauty and Personal Care again displays the lowest accuracy and highest RMSE, indicating it is the hardest category for prediction and ranking.



### Recall@K

- Recall at higher K (e.g., K=50) is extremely high across all categories, but Sports and Outdoors shows the strongest recall across all K values.
- Beauty and Personal Care hits perfect recall at K=50, which is expected in sparse datasets with few items per user.

### NDCG@K

- Electronics leads in NDCG@10 and NDCG@50, meaning its recommendations are ranked most relevantly for top and broader lists.
- Beauty and Personal Care trails in all NDCG values but slightly improves at higher K, likely due to the nature of available item-user relationships.
- Sports and Outdoors shows a strong NDCG@50, confirming robust relevance ranking when expanding the recommendation list.

### MAP@K

- Electronics delivers the highest precision (MAP@10, MAP@20, MAP@50), with Sports and Outdoors close in MAP@50 but trailing in MAP@10 and MAP@20.
- Beauty and Personal Care is lowest in MAP across all cutoffs, indicating it is the most challenging for precise ranking.

**Summary:** Item-based recommenders work best for Electronics and Sports and Outdoors, delivering high accuracy, recall, and ranking precision in both categories. Beauty and Personal Care consistently underperforms, reflecting underlying data sparsity and segment-specific challenge

### Model-Based collaborative Filtering

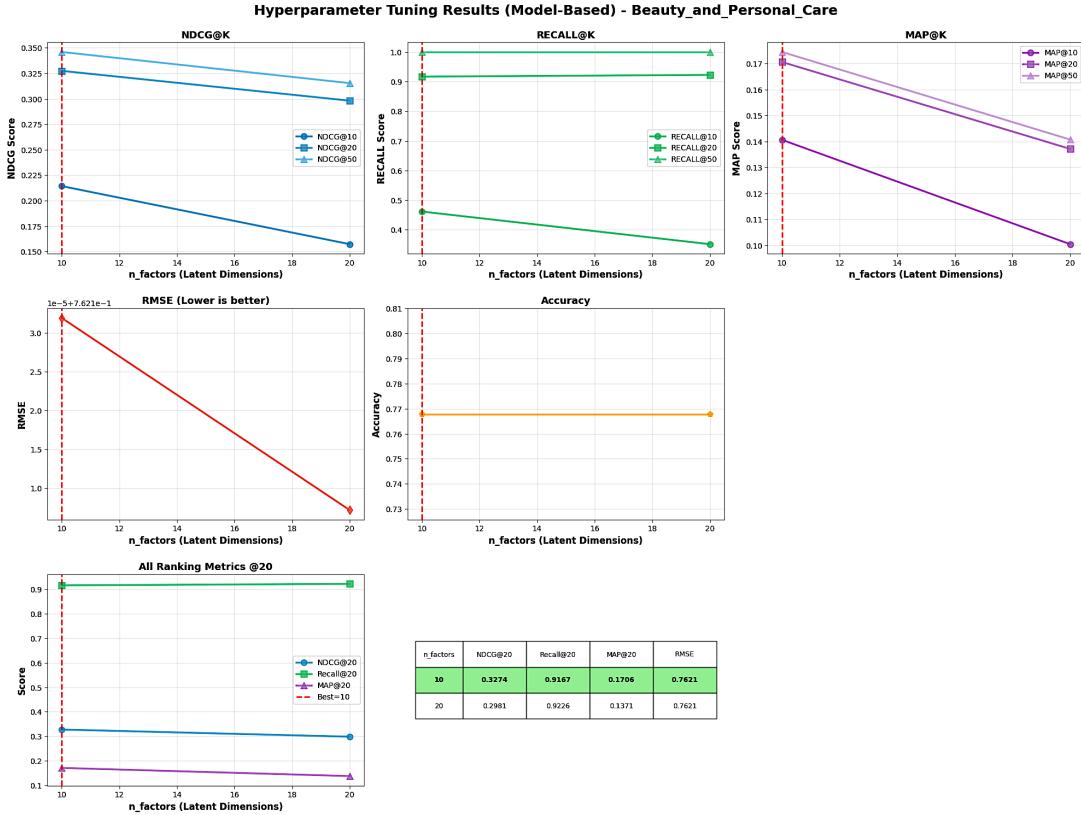
This use SVD Matrix Factorization to decompose rating matrix into latent user and item factors. We used global mean instead of per-user mean for matrix factorization stability. Latent factors (k) is key parameter, indicates tuned on validation set.

```

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-norm

```

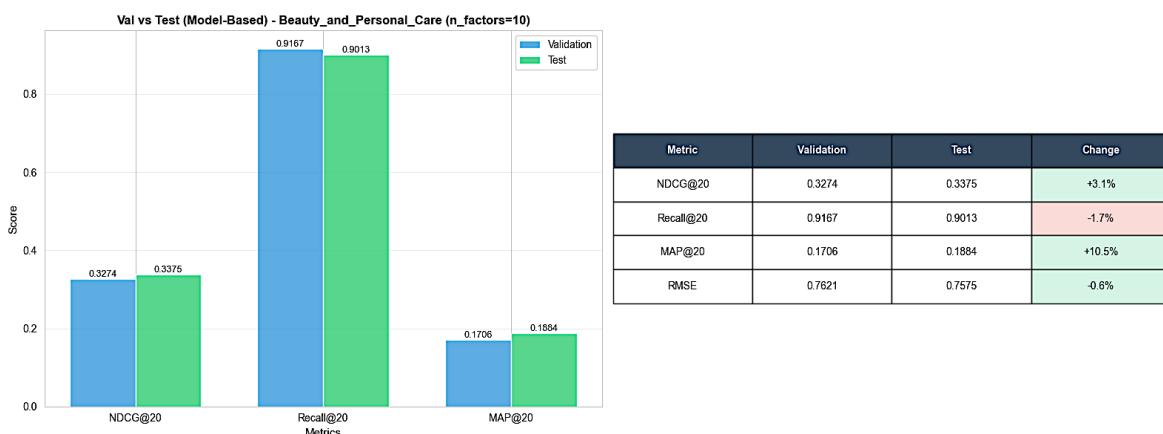
## Result of phase 1: Training and tuning with train dataset



### Key Insights:

- The best overall metrics, NDCG@20 (0.3274), Recall@20 (0.9197), MAP@20 (0.1706), and lowest RMSE (0.7621), are all achieved with the smallest latent dimension setting (10), as highlighted in the summary table.
- Increasing the number of latent factors to 20 leads to a noticeable decline in all ranking metrics, especially MAP@20 (drops from 0.1706 to 0.1371), indicating overfitting or reduced model generalization at higher complexity.
- RMSE remains virtually unchanged between settings, suggesting rating prediction accuracy is insensitive to within this range, but ranking quality suffers with more factors.
- Overall, lower latent dimensions yield better and more stable recommendations for this sparse category, complex models risk degrading top-N recommendation quality.

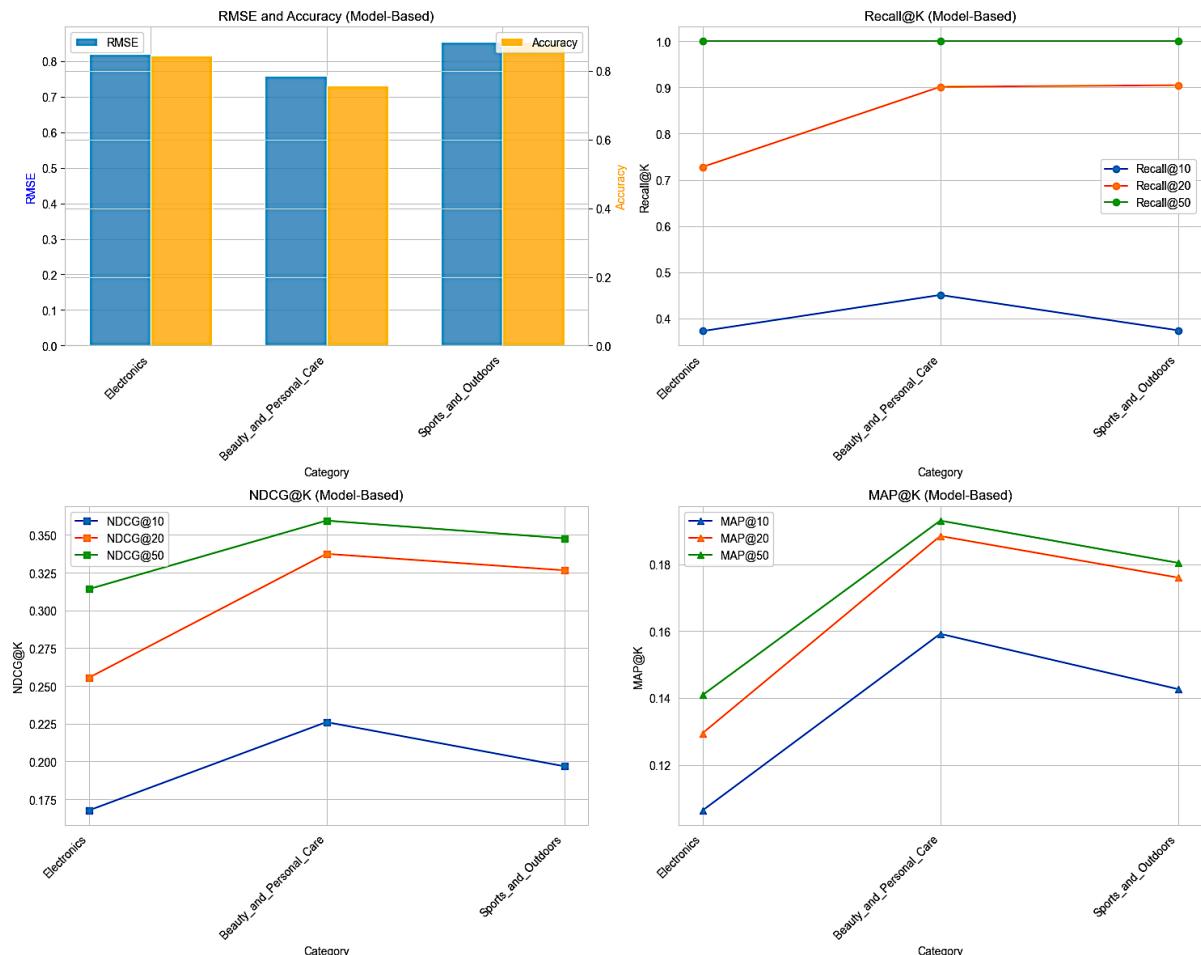
## Result of phase 2: Final evaluation with test dataset



### Key Insights:

- NDCG@20 and MAP@20 both improve on the test set (+3.1% and +10.5%, respectively), indicating the model generalizes well and often ranks truly relevant items higher when evaluated on unseen data.
- Recall@20 decreases slightly (-1.7%) on the test set, implying that while top-N relevance and precision improve, the breadth of captured relevant items shrinks somewhat.
- RMSE drops marginally (-0.6%), showing that prediction accuracy for individual ratings is slightly better on the test set than on validation.
- The gains in ranking metrics (NDCG, MAP) paired with the small drop in recall suggest the SVD model achieves higher precision and more relevant list ordering for the test data, even if fewer items are recovered overall.
- This pattern may result from the model being well-regularized, robust to overfitting, and perhaps benefiting from slight distribution differences between the validation and test split in this dataset

### Result for all categories



### Key Insights:

*Accuracy and RMSE*

- Sports and Outdoors achieves the highest accuracy and lowest RMSE, indicating superior rating prediction and fit for this category.
- Beauty and Personal Care lands at the bottom for accuracy and has the highest RMSE, continuing to be the most challenging domain for both rating and ranking tasks.
- Electronics performs moderately in both accuracy and RMSE metrics.

#### *Recall@K*

- Beauty and Personal Care attains perfect recall at higher K values (), likely related to its data structure (few items per user, easier coverage in top-N).
- Sports and Outdoors shows high recall overall, marginally trailing Beauty and Personal Care at larger K.
- Electronics displays noticeably lower recall, consistent with sparser user engagement or a broader catalog.

#### *NDCG@K*

- Beauty and Personal Care leads on all NDCG metrics, its recommendations are best ranked in terms of true relevance, especially with smaller and moderate N lists.
- Sports and Outdoors ranks close in NDCG, followed by Electronics.

#### *MAP@K*

- Beauty and Personal Care provides the highest precision (MAP@10, MAP@20, MAP@50), especially for top-ranked lists, showing the model is well-calibrated for highly relevant recommendations within this segment.
- Sports and Outdoors has slightly lower MAP than Beauty and Personal Care but still outperforms Electronics, particularly at higher K.
- MAP scores for Electronics lag consistently, reinforcing lower overall recommendation precision for this category.

**Summary:** Model-based recommenders are especially effective for Beauty and Personal Care when ranking and precision matter, but less so for general rating prediction, where Sports and Outdoors excels. Electronics is consistently outperformed by the other categories in both overall accuracy and relevance metrics.

### **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:

$$\text{trending\_score} = \log(\text{rating\_count}) * \text{avg\_rating} * \text{recency\_weight}$$

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 →  $\log(15) \approx 2.7$ , 15K ratings →  $\log(15000) \approx 9.6$  (only 3.6x difference)

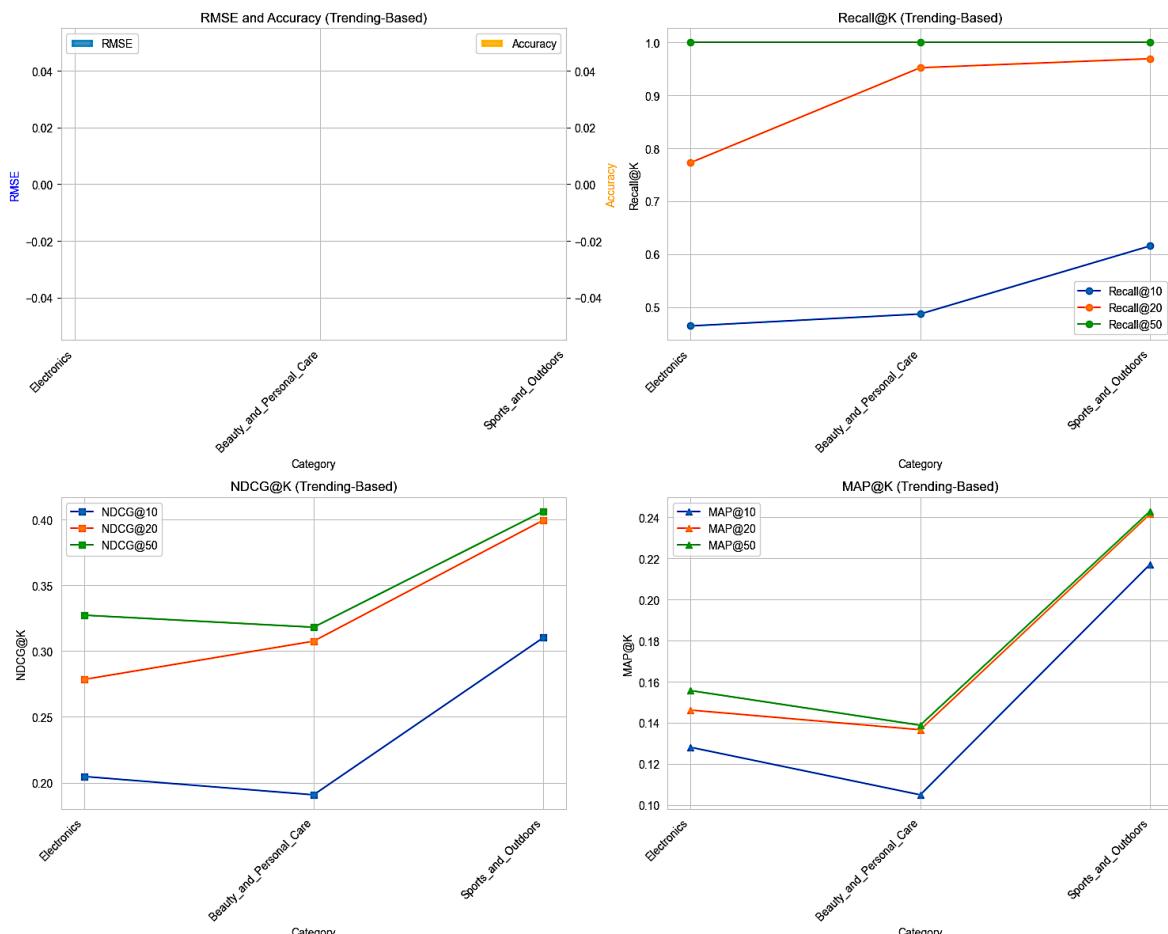
- Allows moderately popular items to compete with blockbusters
- **avg\_rating** - Average rating quality (1-5 scale)
  - Ensures high-quality items rank above low-quality items with similar volume
  - Example: Item A (1000 ratings, 4.8 score) beats Item B (1000 ratings, 3.2 score)
  - Acts as quality filter: popular but poorly-rated items get penalized
- **recency\_weight** - Temporal relevance boost
 

recent\_count = ratings in last 90 days

$$\text{recency\_weight} = 1.0 + 0.5 * (\text{recent\_count} / \text{total\_count})$$
  - 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

### Result for all categories

The trending-based method does not have results for phase 1 training + tuning or phase 2 final evaluation by category because it does not rely on hyperparameters or learning from historical interactions, unlike user/item/model-based recommenders, it simply ranks items based on recent popularity or activity, which can be directly computed for all categories without a dedicated optimization process. Therefore, trending-based results are always “final” and comparable across categories with no distinct training/tuning phases needed.



### Key Insights:

#### *Recall@K*

- Sports and Outdoors leads with highest recall across all K, reaching perfect recall at K=50; users are reliably recommended nearly all relevant trending items in this category.
- Beauty and Personal Care is slightly behind Sports and Outdoors for recall at higher K, still achieving excellent coverage but indicating some limitation in how broadly trends apply to this segment.
- Electronics consistently reports the lowest recall at all K values, showing that trending algorithms may struggle to surface relevant items in a more diverse or less trend-driven domain.

#### *NDCG@K*

- Sports and Outdoors again dominates on ranking metrics (NDCG@10, NDCG@20, NDCG@50), showing high relevance of trending recommendations for top-N lists in this category.
- Beauty and Personal Care shows moderate ranking quality, with scores clustering but not peaking for top-N, indicating solid but unspectacular relevance.
- Electronics lags in NDCG, reinforcing weaker ordering and relevance for recommendations generated by pure trending signals.

#### *MAP@K*

- Trending recommendations are most precise in Sports and Outdoors (highest MAP@K), reflecting a strong alignment of popularity with user interests in this domain.
- Precision (MAP@10, MAP@20, MAP@50) is lowest for Beauty and Personal Care, indicating top-ranking trending items are less likely to match actual user interests or needs.
- Electronics is in the middle, but closer to Beauty and Personal Care than Sports and Outdoors for MAP scores.

**Summary:** Trending-based recommenders are most effective in Sports and Outdoors, achieving high coverage, strong relevance, and best precision overall. Beauty and Personal Care fares adequately in recall and ranking but lacks top precision, while Electronics underperforms, showing that trending signals alone are less suited to heterogeneous or broad product categories.

### **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

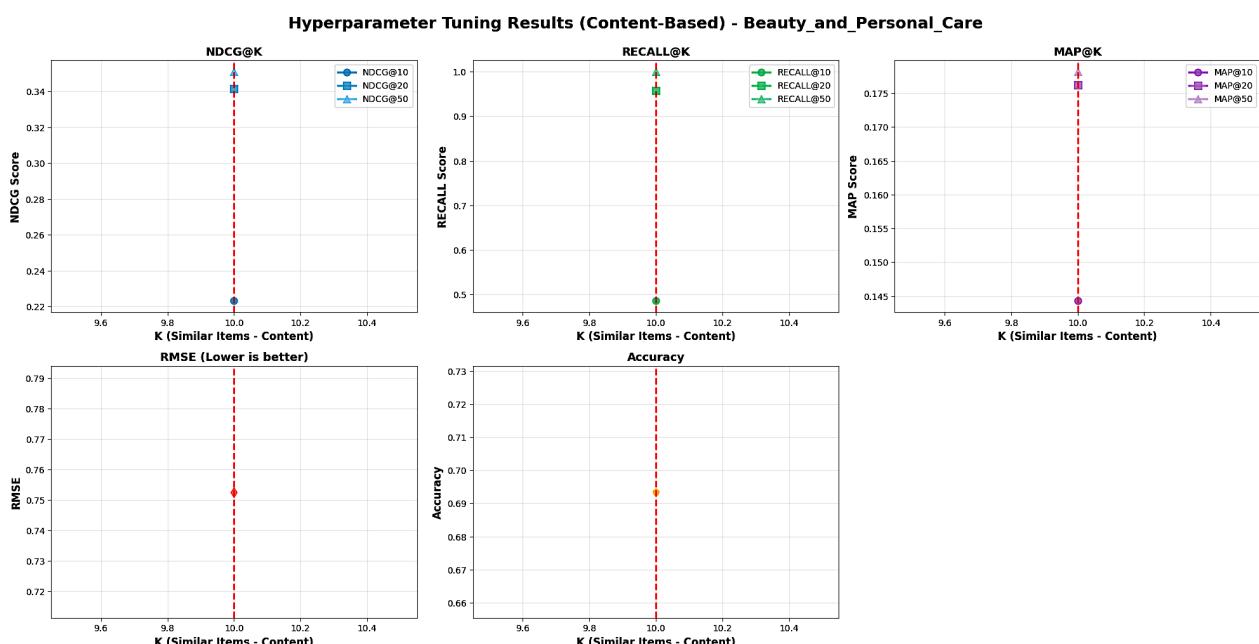
Why TF-IDF over alternatives:

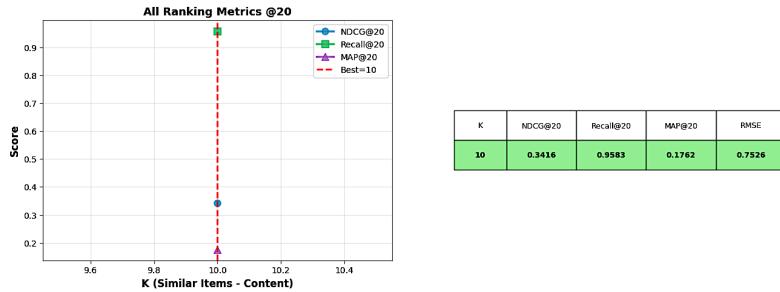
- vs Word2Vec/BERT: Simpler, faster, no pre-training needed
- vs Count Vectorizer: TF-IDF downweights common terms across items
- vs Manual features: Automatically learns important terms from data

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"

Result of phase 1: Training and tuning with train dataset

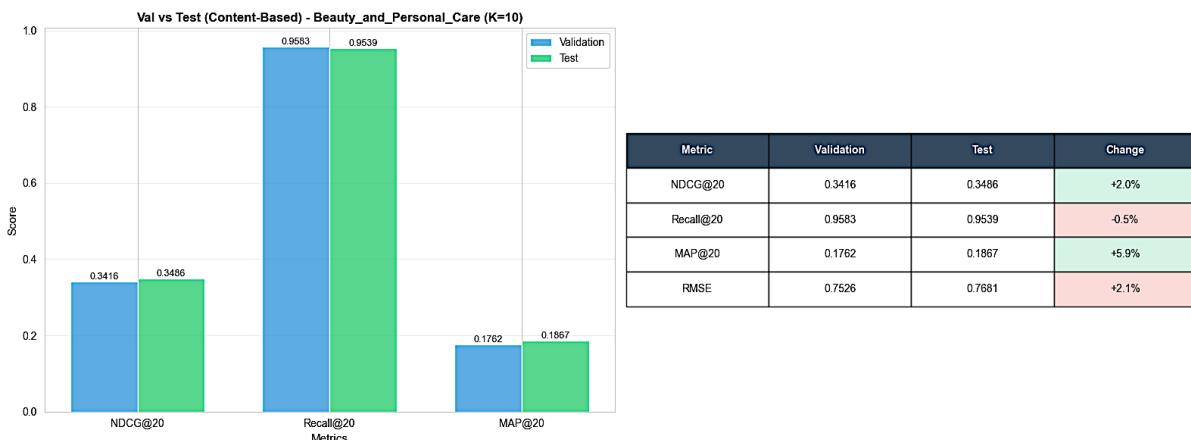




### Key Insights:

- All performance metrics (NDCG, Recall, MAP, RMSE, Accuracy) are reported only for ; there is no observed change with other values, indicating results are either invariant to K or only one setting was considered in this analysis.
- NDCG@20 reaches 0.3416 and MAP@20 is 0.1762, representing strong ranking and precision for recommended items in this category.
- Recall@20 is exceptionally high at 0.9583, showing that content-based methods effectively surface nearly all relevant items for users.
- RMSE (0.7526) and Accuracy (0.6925) reflect competitive rating prediction quality compared to user/item/model-based approaches, reinforcing the strength of leveraging item features in sparse settings.
- The lack of a tuning curve suggests either content-based similarity yields stable results regardless of K, or K is not a meaningful hyperparameter for this method given the feature setup here.

### Result of phase 2: Final evaluation with test dataset

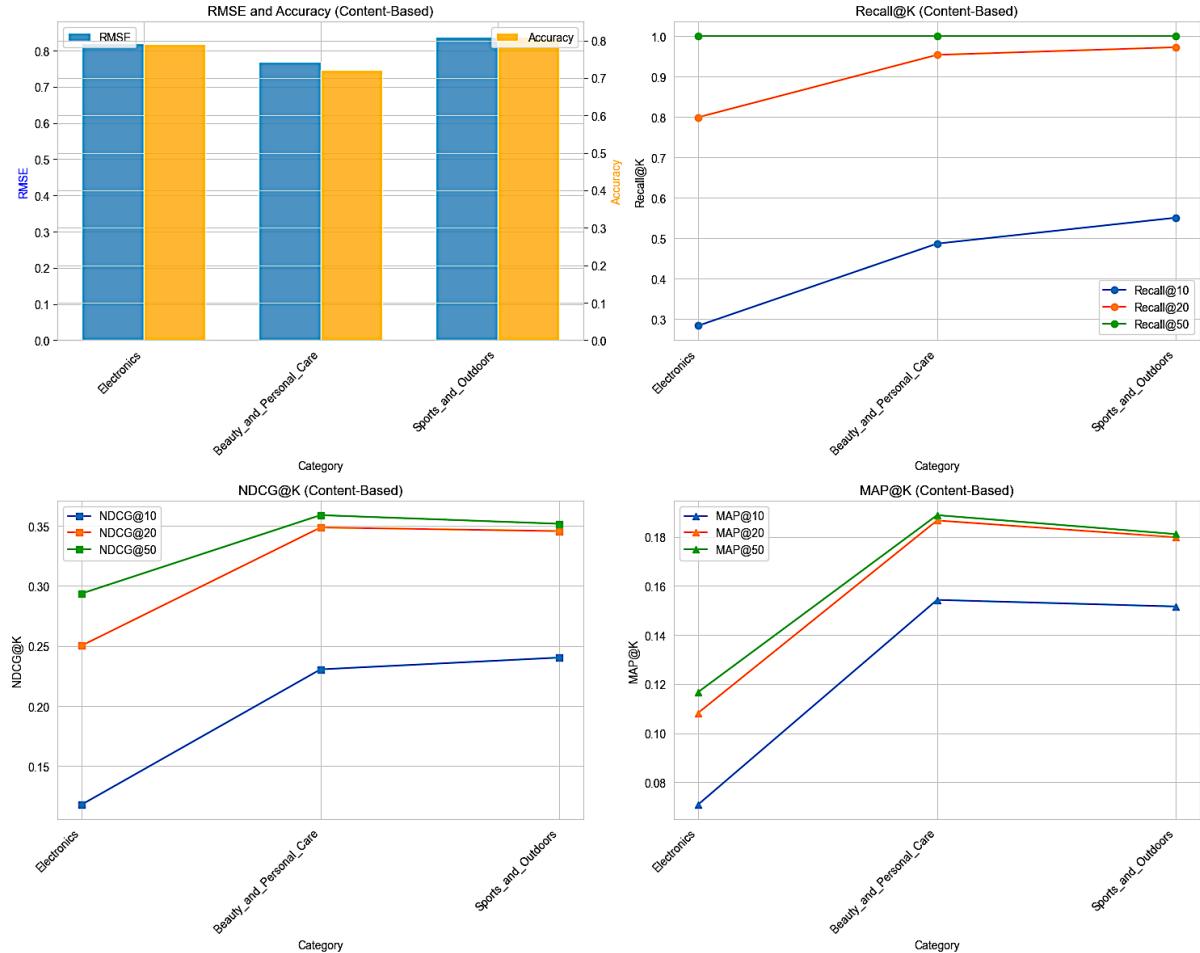


### Key Insights:

- Both ranking quality (NDCG@20) and precision (MAP@20) improve on the test set compared to validation, with NDCG@20 rising by 2% and MAP@20 by nearly 6%, suggesting the model may generalize well beyond the tuning data.
- Recall@20 remains very high and decreases only slightly (-0.5%), which means the vast majority of relevant items are still correctly retrieved despite the data split.
- RMSE increases by 2.1% on the test set, showing a marginal reduction in prediction accuracy for rating scores but still within a competitive range for a content-based approach.

- The overall improvements in NDCG and MAP suggest that the content-based method, by leveraging item features, is not overly susceptible to overfitting and may even benefit from mild distributional differences between validation and test set.
- These patterns highlight the effectiveness and robustness of content-based approaches, especially in sparse or cold-start-heavy environments found in Beauty and Personal Care recommendations.

### Result for all categories



### Key Insights:

#### *RMSE and Accuracy*

- Sports and Outdoors and Electronics achieve higher accuracy (above 0.8) and lower RMSE than Beauty and Personal Care, indicating stronger rating prediction capabilities in these domains.
- Beauty and Personal Care continues to lag with the lowest accuracy and the highest RMSE, a trend observed in other recommendation families as well.

#### *Recall@K*

- Sports and Outdoors and Beauty and Personal Care both achieve near-perfect recall at high K, reflecting strong retrieval of relevant items for most users.
- Electronics, though robust in ranking and precision, has much lower recall, indicating a challenge to cover all relevant items for users in broader/non-niche categories.

### *NDCG@K*

- Beauty and Personal Care shows the top NDCG@10, NDCG@20, and NDCG@50 scores, making it the best category for ranking relevance with content-based recommenders.
- Sports and Outdoors closely follows, while Electronics often trails in NDCG scores, suggesting less optimal ordering despite reasonable accuracy.

### *MAP@K*

- Precision at every cutoff (MAP@10, MAP@20, MAP@50) is highest for Beauty and Personal Care and Sports and Outdoors, confirming these segments benefit most from content-driven similarity and feature-based personalization.
- Electronics consistently achieves lower MAP scores, highlighting its challenge for precise top-N recommendations in a highly diverse or frequently changing product set.

**Summary:** Content-based recommenders excel most in Beauty and Personal Care and Sports and Outdoors, delivering high recall, relevance, and precision; these approaches face greater difficulty in Electronics, where user preferences and item features might be more fragmented or less predictive.

### *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

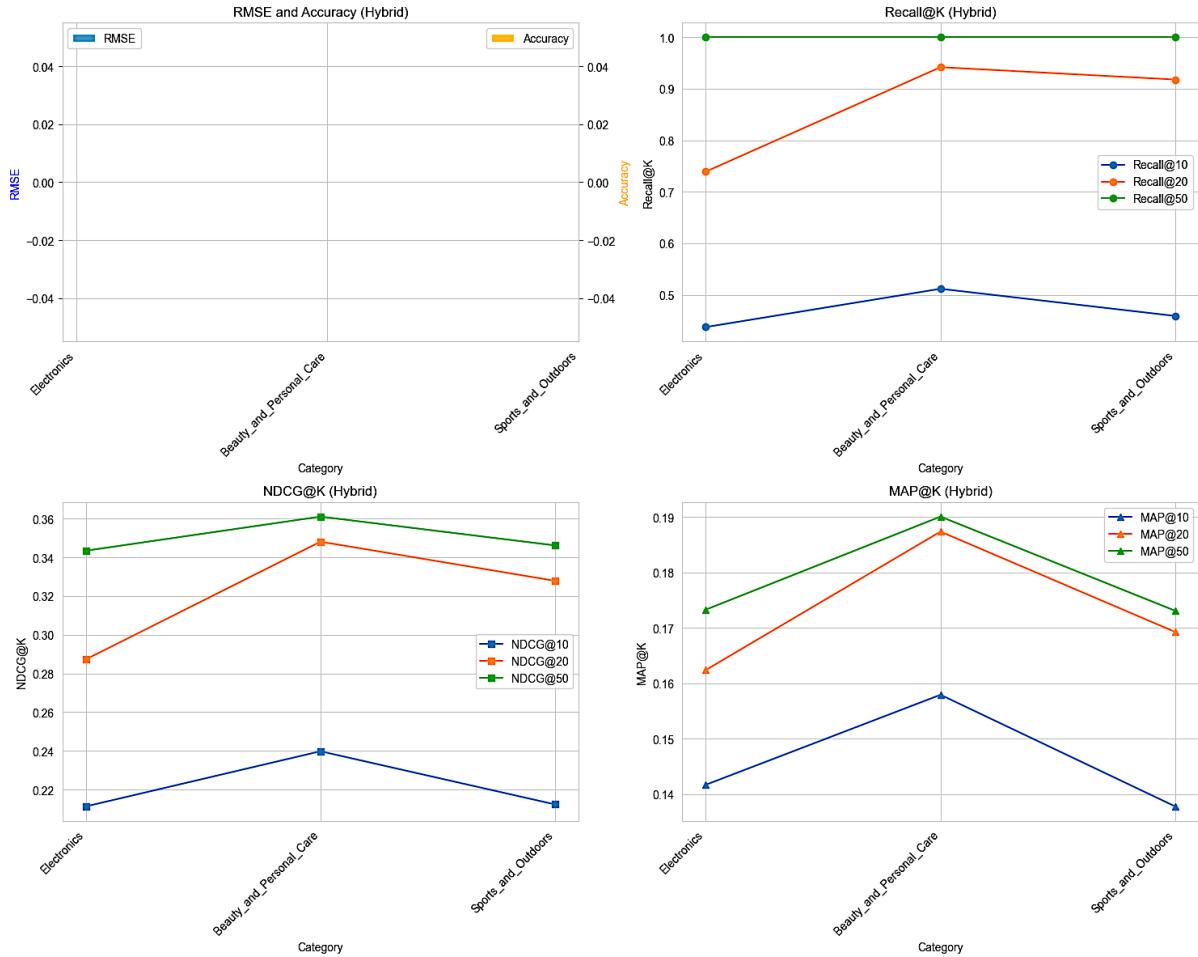
```
# Detect scenario
scenario = detect_scenario(user, threshold=5)

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

### Results for all categories:



### Key Insights:

#### *Recall@K*

- Both Sports and Outdoors and Beauty and Personal Care achieve extremely high recall at broader cutoffs (Recall@50 nearly 1.0), indicating these hybrid recommenders can successfully retrieve almost all relevant items for most users at larger recommendation list sizes.
- At smaller K (Recall@10, Recall@20), Beauty and Personal Care stands out for achieving higher recall compared to others, showcasing strong top-N coverage in this segment.
- Electronics consistently trails in recall across all K, reflecting persistent challenges for broad relevant item recovery in this category.

#### *NDCG@K*

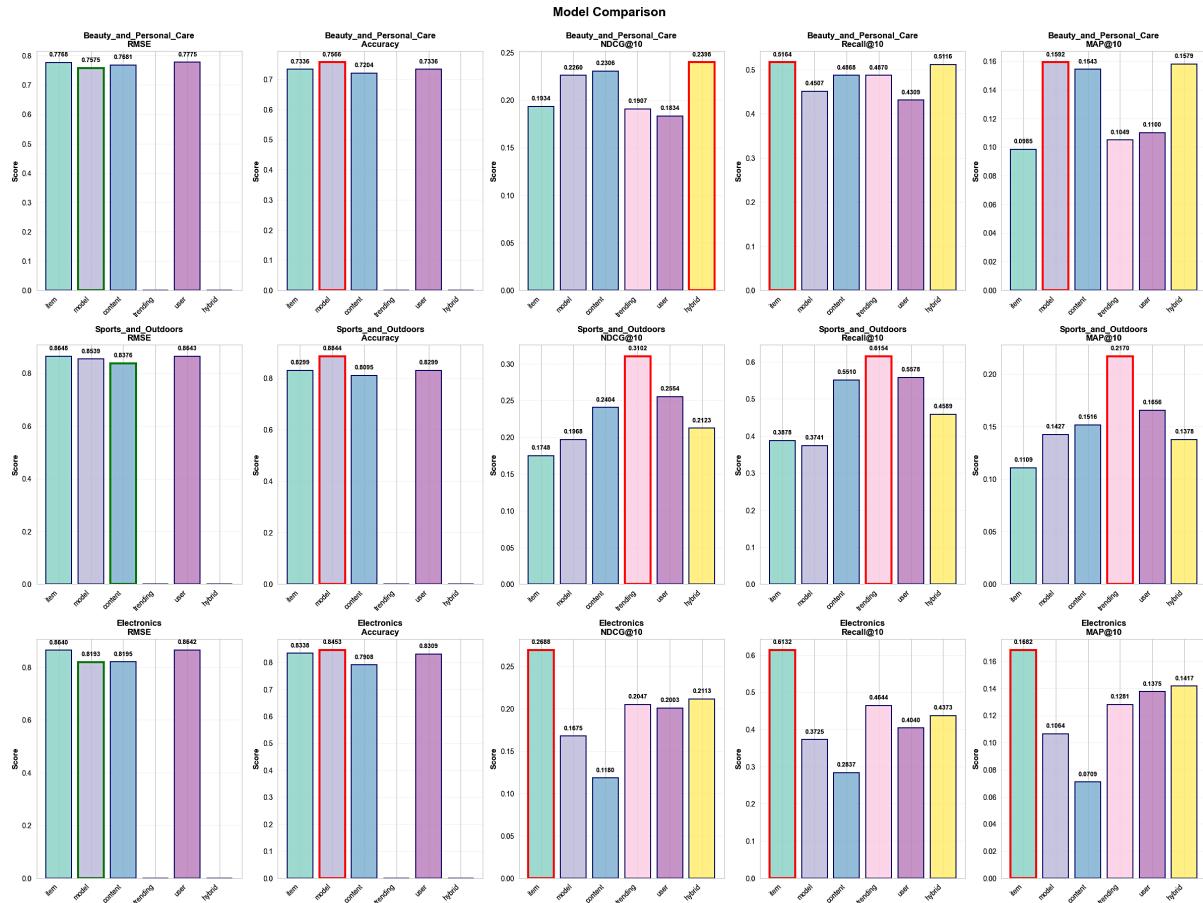
- Beauty and Personal Care achieves the highest NDCG@20 and NDCG@50, solidifying its status as the most robust environment for hybrid ranking relevance.
- Sports and Outdoors remains competitive and shows only slight declines compared to Beauty and Personal Care; Electronics lags, especially at lower K cutoffs.

#### *MAP@K*

- Precision (MAP@10, MAP@20, MAP@50) is again highest in Beauty and Personal Care, confirming the hybrid method's strong alignment between feature-based and collaborative signals for what is relevant.
- Sports and Outdoors delivers nearly comparable MAP at broader cutoffs, making it an advantageous environment for the hybrid model, particularly at K=20 and K=50.
- Electronics has the lowest MAP in every scenario, showing less precise and relevant recommendations in this more general category.

**Summary:** The hybrid recommender leverages the best of both worlds, strong collaborative and content signals, and this pays off especially in Beauty and Personal Care (highest recall and precision) and Sports and Outdoors, while Electronics remains challenging due to its heterogeneous and potentially less-structured item space.

### Compare among algorithms



### Key Differences Between Methods

- Trending-based models stand out for Recall@10 and Recall@50, being best in Electronics and Sports and Outdoors for finding popular items, but lag behind in ranking quality (NDCG/MAP) and accuracy for individual preference nuances.
- Model-based (SVD) consistently achieve the highest accuracy in Beauty and Personal Care but are not always the top in MAP or NDCG@10, indicating their strength is rating prediction more than top-N recommendation.

- Item-based and Content-based methods compete closely in most ranking and precision metrics, with Content-based especially powerful in Beauty and Personal Care and Sports and Outdoors, reflecting high feature relevance in these domains.
- Hybrid models generally do not dominate any one metric, but show consistently good performance across the board, often landing just shy of the best result for each evaluation, thus offering balanced accuracy and ranking without major weaknesses, while handling well cold-start problem.

#### Cross-category Patterns

- Beauty and Personal Care often shows the largest spread between methods, with Content-based and Model-based solutions excelling; Item-based and Trending models are less suited here.
- Sports and Outdoors displays strong performance for Trending and Content-based methods, likely due to fewer, highly-engaged items and users, making these approaches especially effective.
- Electronics is outperformed by other categories in ranking and precision, regardless of the algorithm, indicating intrinsic challenges like wider item diversity or user preference variance.

**Summary:** No single method dominates everywhere, each model shines in specific metrics or for specific categories due to data structure, user behavior, and item features. Content-based and model-based methods thrive in Beauty and Personal Care, trending-based and content-based are especially effective in Sports and Outdoors, and Electronics remains a broadly challenging domain for all approaches.

Finally, all models for each algorithm, 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...

```

### 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.

```
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
```

## 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.

```
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
```

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)

---

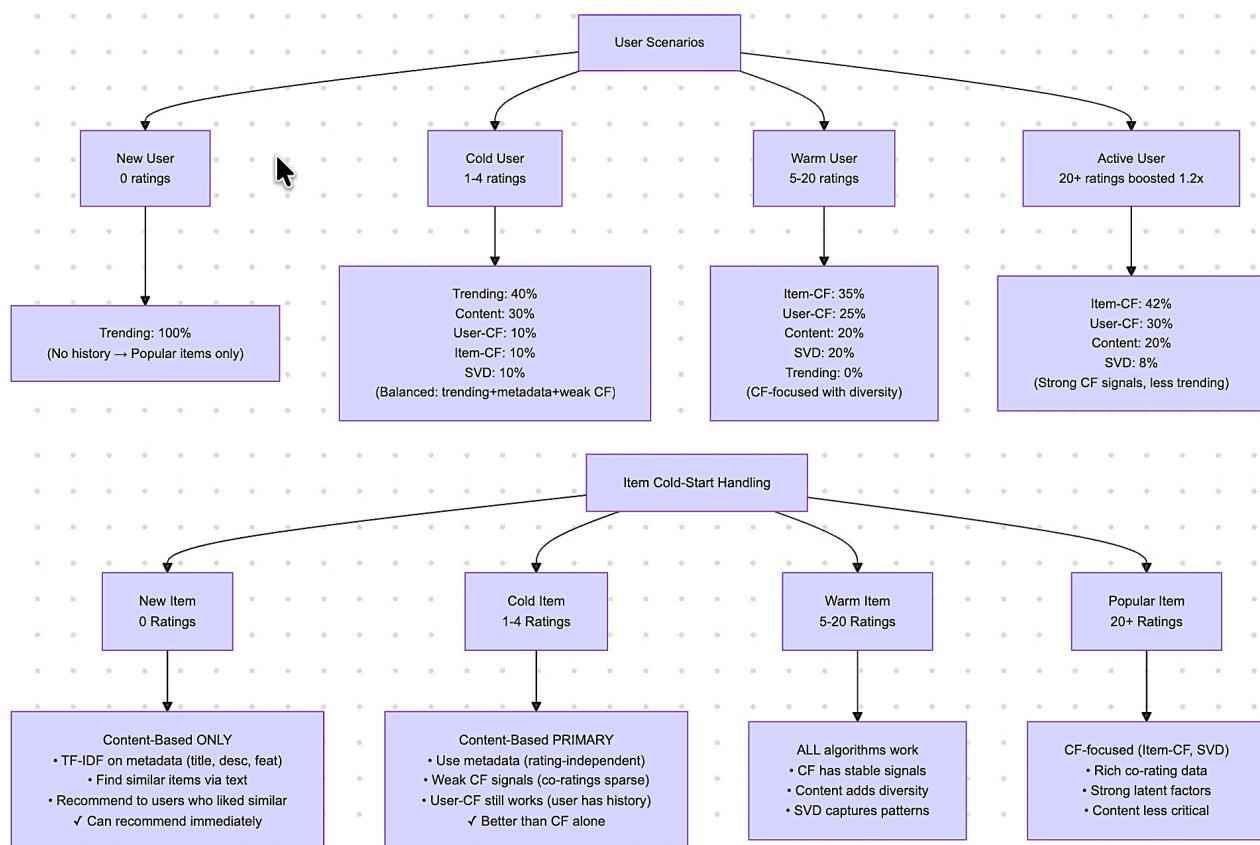
##### USER SCENARIOS:

- 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

The screenshot shows the homepage of the Amazon Product Recommendation System. At the top, there's a navigation bar with categories like Electronics, Sports & Outdoors, and Beauty & Personal Care. Below the navigation is a grid of five product cards, each featuring a sponsored product and its details (e.g., price, rating, and 'Add to Cart' button). In the footer, there's a 'Get to Know Us' section, a 'Recommendation System' section (which lists User-Based CF, Item-Based CF, Content-Based, Model-Based, Trending-Based, and Hybrid Model), a 'Categories' section (listing Electronics, Sports & Outdoors, and Beauty & Personal Care), and a 'Project Info' section (listing Northeastern University, IE7276 Data Mining - Fall 2025, Group 6, and Quoc Hung Le & Matthew Eckert). The 'Recommendation System' and 'Project Info' sections are highlighted with red boxes.

### Create or Login

The screenshot shows the login page of the Amazon Product Recommendation System. A modal window titled 'Create account' is open, prompting the user to enter their name ('le.quo'), email ('le.quo@northeastern.edu'), and password ('\*\*\*\*'). A large red box highlights the input fields for name, email, and password. Below the form, there's a note about agreeing to Amazon's Conditions of Use and Privacy Notice, and links for 'New to Amazon?' and 'Already have an account? Sign in'. The background shows a blurred view of the product recommendation grid, with one item from the previous screenshot (the Transcend TS-RDF5K USB 3.1 SDHC/SDXC/microSDHC/SDXC Card Reader) clearly visible.

**When no user login, use trending**

## New user, hybrid with trending 100%

amazon.ml

Deliver to Boston 02118

Beauty & Personal Care Search Amazon

Hello, le quo  
Account & Lists

Returns & Orders Sign Out

All Today's Deals Customer Service Registry Gift Cards Show Cold Items Only Electronics Sports & Outdoors Beauty & Personal Care

New User  
No rating history - showing popular trending items

Algorithm Strategy: hybrid-new-user-trending(100%)  
Combining multiple algorithms based on your profile

Recommended for you in Beauty & Personal Care

Trending THAYERS Alcohol-Free, Hydrating Coconut Water Witch Hazel Facial Toner with Aloe Vera Formula, 1... \$10.95 prime FREE delivery 100% Match

Trending ILNP Playdate - Vivid Magenta Holographic Jelly Nail Polish \$10.00 prime FREE delivery 33% Match

Trending APRS sponsored L'Oréal Paris Buildable Voluminous Original Volume Building Mascara, Black Brown, 0.28 fl. oz. \$6.35 prime FREE delivery 25% Match

Trending APRS sponsored Maybelline Instant Age Rewind Eraser Dark Circles Treatment Concealer, Fair 0.2 oz (Pack of 2) \$7.73 Price not available prime FREE delivery 20% Match

Trending APRS sponsored essence Iash Princess False Lash Effect Mascara | Gluten & Cruelty Free \$4.99 prime FREE delivery 17% Match

## Active user, adaptive hybrid approach

The screenshot shows the Amazon.mil website interface. At the top, there's a navigation bar with 'Electronics' selected. A red box highlights the 'Account & Lists' button and the '>Returns & Orders' link. Below the navigation, a message says 'Hello, AHL3STEVTNSEZI7TWZRLAO4EULKKA'. A red box highlights the 'Active User' badge with the text '25 ratings - highly personalized'. Another red box highlights the 'Algorithm Strategy: hybrid-warm-user-item(51%)+content(24%)+model(24%)' message. The main content area is titled 'Recommended for you in Electronics'. It displays five product cards, each with a 'For You' badge. The products include a camera lens protection filter, a network switch, a Western Digital hard drive, a USB hub, and a Fire HD tablet.

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

This comparison shows two versions of an Amazon search results page for 'Electronics'. The left side, labeled 'Before', shows a network switch in a 'Trending' section. The right side, labeled 'After', shows the same network switch moved to a 'For You' section after the user rated it. Other products visible include SSDs, a tablet, and a portable hard drive.

## Cold-item handling

Click to tick for 'Show Cold items Only'

**Cold Items - By Training Rating Count**  
Items grouped by number of ratings in training set. Scroll right to see more items in each row.

**4 Ratings** (page 1 of 2 - showing 5 of 10 items)

- APRS sponsored  
Maplefour Water Bottle Caps, Water Jug Caps, 24 Pack 55mm 3 & 5 Gallon Water Bottles...  
★★★★ 4.4  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**
- APRS sponsored  
FIFTY 8/12 Pounds Multifunctional Hand and Forearm Trainer - Spinning Burn Muscle Trainin...  
★★★★ 4.6  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**
- APRS sponsored  
The Original Glow in The Dark Smiley Face PVC Rubber Morale Patch by NEO Tactical Gear  
★★★★ 4.7  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**
- APRS sponsored  
Skyweo Exercise Band Set, Resistance Bands Set Long Latex Elastic Bands Wide Fitness...  
★★★★ 4.3  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**
- APRS sponsored  
MOSHAY Roll Over Image to Zoom in Bicycle Training Wheels Fits 16 18 20 22 24 inch Kids...  
★★★★ 3.8  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**

**3 Ratings** (page 1 of 2 - showing 5 of 10 items)

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

**Cold Items - By Training Rating Count**  
Items grouped by number of ratings in training set. Scroll right to see more items in each row.

**4 Ratings** (page 1 of 4 - showing 5 of 20 items)

- APRS sponsored  
Easton Team Player Bag  
★★★★ 4.4  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**
- APRS sponsored  
Disney Frozen Olaf Jumbo Tritan Hydration Bottle (739ml)  
★★★★ 4.1  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**
- APRS sponsored  
Gould & Goodrich 840 Gold Line Handcuff Case With Belt Loop [Chestnut Brown] Holds most...  
★★★★ 4.6  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**
- APRS sponsored  
Deep Blue Gear Kids Ultra Dry 2 Snorkel, Junior  
★★★★ 3.7  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**
- APRS sponsored  
PS Products Undercover Concealment Holster with Adjustable Neck Cord, Black  
★★★★ 4.6  
Price not available  
**prime** FREE delivery  
**Rate Product** **Add to Cart**

**3 Ratings** (page 1 of 4 - showing 5 of 20 items)

## 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

## 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