# IRWA Project Part 2: Indexing and Evaluation Report

Authors: Lucas Andreu, Pau Chaves, Pol Bonet, Joan Company

Course: Information Retrieval and Web Analytics (IRWA)

Date: November 2, 2025

## 1. Project Summary

This second part extends the previous preprocessing work by implementing the indexing, ranking, and evaluation modules of our fashion-product search engine. The system builds an inverted index using conjunctive (AND) queries, applies a TF-IDF ranking model, and evaluates retrieval quality using standard IR metrics.

## 2. Indexing

### 2.1 Inverted Index Construction

After preprocessing the 28 080 documents, we built a Boolean inverted index where each term ID maps to the list of document IDs that contain it. Only documents that contain all query terms are considered (AND logic).

| Metric | Value |
|---|---|
| Documents indexed | 28 080 |
| Vocabulary size | 9 048 |
| Indexed fields | title, brand, category, description, metadata |
| Saved files | boolean_inverted_index.json, docid_pid_map.json, indexed_fields.json |

## 2.2 Proposed Test Queries

We created five realistic test queries using an automated query-selection technique that first defines limited attribute lexicons (gender, category/type, color, material, fit, style), then normalizing and filtering applicants whose stemmed tokens from the index's same preprocessing are absent, (3) removing any candidates whose stemmed tokens from the index are not present in the (4) scoring remaining candidates based on the sum of document frequencies (DF) as a rapid popularity indicator; inverted vocabulary; and (5) creating queries from attribute templates. (e.g. gender + type + attribute) filled with high-DF terms. Every made-up query is checked to make sure all the words are in the dictionary and that a Boolean requirement is met. AND search provides at least one result, therefore every suggested query produces actual results; the generator enforces diversity (limits recurring main types), uses a limited If necessary, set of hand fallbacks; prints tokens/DF/approximate matches for inspection; saves the last list to proposed_test_queries.json for reproducibility:

| # | Query | Estimated Matches |
|---|-------|-------------------|
| 1 | men shirt regular blue | ≈ 50 |
| 2 | men shirt solid blue | ≈ 50 |
| 3 | women cotton kurta straight | ≈ 50 |
| 4 | men slim fit formal shirt | ≈ 50 |
| 5 | men running shoes black | ≈ 28 |

Bias note. DF-guided selection tends to favor popular attributes; we mitigate this by mixing style/material/color terms and varying lengths (4–5 tokens).

## 2.3 TF-IDF Ranking Implementation

A log-scaled TF-IDF model with cosine similarity was implemented. Top results for two validation queries ('women full sleeve sweatshirt cotton' and 'men slim jeans blue') were semantically correct, validating the ranking logic.

Query: women full sleeve sweatshirt cotton

Top result: Full Sleeve Solid Women Sweatshirt  (score = 0.907)

Query: men slim jeans blue

Top result: Slim Men Blue Jeans  (score ≈ 0.717)

These results confirm the correctness of tokenization, weighting, and cosine similarity computation. Ranked lists are stored in ranked_results.json.

$$tf_{i,j} = 1 + \log_2(f_{i,j}), \quad idf_i = \log_2\left(\frac{N}{df_i}\right)$$

# 3. Evaluation

## 3.1 Implemented Metrics

Seven evaluation metrics were implemented: Precision@K, Recall@K, Average Precision@K, F1@K, MAP, MRR, and NDCG@K. Each metric captures a distinct performance aspect.

Metrics were computed at K = 20 for validation queries and K = 10 for expert queries.

## 3.2 Results for Validation Queries (K=20)

| Query | P@20 | R@20 | AP@20 | F1@20 | NDCG@20 | MRR |
|---|---|---|---|---|---|---|
| Q1 – women full sleeve sweatshirt cotton | 0.65 | 1.00 | 0.694 | 0.788 | 0.873 | 1.00 |

| Q2 – men slim jeans blue | 0.50 | 1.00 | 0.627 | 0.667 | 0.833 | 1.00 |
| --- | --- | --- | --- | --- | --- | --- |

Summary: MAP = 0.66, MRR = 1.00, mean NDCG@20 = 0.853

Both validation queries achieve perfect recall and ideal first-hit rank, proving the correctness of the index and ranking implementation.

## 3.3 Expert-Labeled Queries (K=10)

| QID | P@10 | R@10 | AP@10 | F1@10 | NDCG@10 | MRR |
| --- | --- | --- | --- | --- | --- | --- |
| 1 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 2 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| 3 | 0.40 | 1.00 | 0.43 | 0.571 | 0.588 | 0.25 |
| 4 | 0.80 | 1.00 | 0.97 | 0.889 | 0.992 | 1.00 |
| 5 | 0.10 | 1.00 | 0.20 | 0.182 | 0.387 | 0.20 |

Summary: MAP = 0.32, MRR = 0.29, mean NDCG@10 = 0.393

Queries 3 and 4 display excellent ranking behavior, with near-perfect ordering (NDCG $\approx$ 1). Queries 1 and 2 have no relevant labels, lowering the aggregate results.

### Manual Labeling Procedure

To evaluate retrieval quality under realistic conditions, we manually assigned binary relevance labels (1 = relevant, 0 = non-relevant) for the top-10 results of each expert query. Labeling decisions were based on the match between the query intent and product attributes such as type, gender, material, and category. For example, a query "men slim fit formal shirt" was considered relevant only if the product was a men's formal shirt with a slim fit, excluding similar but mismatched items (e.g., women's shirts or casual wear). This manual assessment provided a reliable ground truth for evaluating ranking metrics.

## 3.4 Ablation Results

| Configuration | MAP | MRR | mean NDCG@10 |
|---|---|---|---|
| tfidf_and | 0.600 | 0.600 | 0.600 |
| tfidf_no_meta | 0.483 | 0.600 | 0.565 |
| tfidf_title3 | 0.326 | 0.290 | 0.395 |
| tfidf | 0.320 | 0.290 | 0.393 |

The ablation confirms that query formulation and field selection have a greater impact than minor weighting changes.

## 3.4 Final Analysis and Conclusions

The evaluation phase allowed us to assess the effectiveness of our retrieval models using both predefined and manually labeled queries. Across all experiments, we observed consistent behavior among ranking metrics (MAP, MRR, and NDCG), confirming the stability and reliability of the retrieval system.

In the baseline TF-IDF configuration, the system achieved moderate effectiveness (MAP = 0.32, MRR = 0.29, NDCG@10 = 0.39). Subsequent ablation experiments revealed substantial performance differences depending on the retrieval strategy and field weighting.

The Boolean AND filtered TF-IDF configuration achieved the best overall results (MAP = 0.60, MRR = 0.60, NDCG@10 = 0.60), demonstrating that restricting candidate documents to those containing all query terms significantly improves retrieval precision and ranking quality.

Excluding metadata fields (tfidf_no_meta) also produced competitive results (MAP = 0.48), suggesting that metadata may introduce lexical noise. Conversely, heavily weighting the title field (tfidf_title3) did not yield substantial gains, indicating that product descriptions remain valuable for contextual discrimination.

Overall, the results confirm that:

- Balanced weighting between title and description fields provides robust performance.
- Applying Boolean AND filtering before ranking improves relevance by reducing false positives.

Metadata fields require careful curation to avoid lowering precision.

These findings highlight that a combination of effective text preprocessing, appropriate field weighting, and term-based filtering leads to a more accurate and interpretable search engine for the fashion-products domain.

## 4. Metric Interpretation & Analysis

Each evaluation metric highlights a different dimension of system performance:

Precision@K shows ranking accuracy. Recall@K demonstrates that the index retrieves all relevant documents. AP and MAP measure average ranking quality across queries. F1 balances precision and recall. MRR focuses on how early users find a relevant result. NDCG evaluates position-weighted ranking correctness.

On the validation set, the system reached perfect recall (1.0) and high NDCG (> 0.85), showing excellent ordering of relevant results.

The expert-labeled queries show similar tendencies but reveal that performance depends strongly on query clarity and label completeness.

## 5. System Limitations & Improvements

| Issue | Impact | Proposed Solution |
|---|---|---|
| Strict AND matching | Misses partially relevant docs | Allow OR logic or BM25 ranking |
| No stemming/lemmatization | Fails to match variants | Add Porter stemmer or lemmatizer |

| Equal field weights | Metadata noise hurts precision | Use field boosting |
| --- | --- | --- |
| No semantic understanding | Misses synonyms | Integrate word embeddings |
| Limited labels | MAP underestimation | Expand relevance judgments |

## 6. Conclusions

The TF-IDF AND-based system satisfies all requirements of Part 2 and produces reliable baseline performance. Validation results: MAP = 0.66, NDCG@20 = 0.85, Recall = 1.0 → the model retrieves and ranks relevant documents correctly. Expert results: MAP = 0.32 with some variability due to missing labels, but excellent results for well-defined queries (e.g., Q4 ≈ 0.99 NDCG). The ablation study highlights that field weighting and candidate selection are the most influential components. Future iterations will incorporate BM25, term normalization, and semantic query expansion to further improve precision and robustness.