

Performance Evaluation in (CS) IR System

performance Evaluation in IR is the process of measuring how well an information retrieval system retrieves relevant documents for a user's query.

It helps us Compare systems, improve algorithms, and ensure users get accurate results.

Purpose of IR Evaluation System

- ✓ To check accuracy of retrieved results.
- ✓ To compare different retrieval models.
- ✓ To improve ranking algorithms
- ✓ To assess User Satisfaction

Key Concepts

- * Relevant documents
 - documents that actually match the user's info need.
- * Retrieved documents
 - documents returned by system (whether relevant OR NOT)

True/False classification

Category

True positive (TP)

False positive (FP)

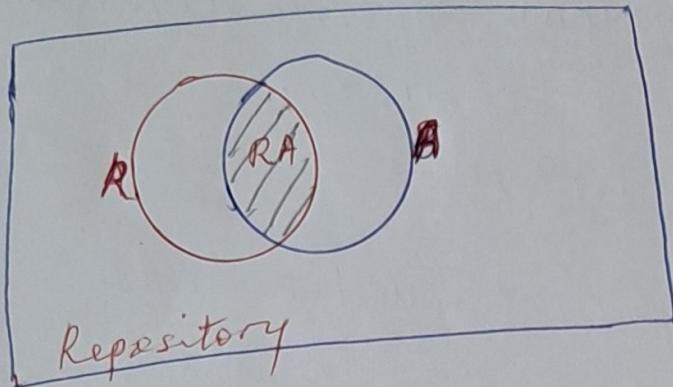
False Negative (FN)

True Negative (TN)

Meaning

- ✓ Relevant documents Retrieved
- ✓ NOT-relevant Retrieved
- ✓ Relevant but NOT Retrieved
- ✓ Non-relevant and not retrieved
(Usually ignored in IR)

	Relevant	NOT-Relevant
Retrieved	TP	FP
NOT-Retrieved	FN	TN



RA = Retrieved Relevant

R = Relevant

A = Retrieved

Core Evaluation Metrics

Precision :- Measures accuracy :- of all retrieved documents, how many are relevant?

$$\text{precision} = \frac{RA}{A} = \frac{TP}{TP + FP}$$

High precision \Rightarrow More relevant results at the top.

Recall :- Measure completeness :- of all relevant documents, how many were retrieved?

$$\text{Recall} = \frac{RA}{R} = \frac{TP}{TP + FN}$$

High Recall \rightarrow System retrieves most of the relevant items.

F1-Score :- Harmonic mean of precision and recall.

$$\text{F1-Score} = \frac{2 \times (\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}}$$

Performance Evaluation in IR System is commonly measured in terms of :-

Efficiency and Effectiveness.

Efficiency:- "Doing things the right way" refers how well the system uses time, Memory and Computation to produce results.

- ✓ Fast Searching (low response time)
 - The system returns result quickly, even for large collections.
 - Efficient indexing and query processing techniques are used.

Example: Using an inverted index instead of scanning all documents

- ✓ High throughput
 - The system can handle many queries per second
 - Efficient algorithms allow many users to search at once.

Example: Google processing millions of queries simultaneously

- ✓ Efficient Use of Memory & Storage
 - Index structures are compact and optimized
 - Unnecessary data is not stored

Example: Using compression for posting lists

✓ Scalable Architecture

- The system grows well with more data & more users
- Performance doesn't collapse when the collection size increases.

Example: Distributed indexing in search engines

✓ Minimizing Computational Cost

- Algorithms avoid heavy operations
- Query processing is optimized

Example: Using Skip pointers to avoid scanning long posting lists

Summary:

* "Doing things the right way" in an IR system means retrieving information efficiently — fast, scalable, and resource friendly — without wasting time, memory, or computing power

Effectiveness = 'Doing the right things'

Effectiveness measures how well the system achieves its intended goals

In IR:

- ✓ Does the system return relevant documents?
- ✓ How accurate are the result?
- ✓ Common Metrics

precision, Recall, F1-score & Accuracy.

* Effectiveness focuses on quality & accuracy

Precision = $\frac{\text{the number of relevant retrieved documents}}{\text{Total number of retrieved documents}}$

Interpretation

"of all documents the system retrieved, how many were actually relevant?"

High precision is the system retrieves mostly relevant documents

low precision is the system retrieves many irrelevant documents (noise)

Example Suppose a system retrieves 10 documents for a query.

Relevant Retrieved = 7

IR Relevant Retrieved = 3

$$\text{So precision} = \frac{7}{10} = 0.7 = 70\% \quad \equiv$$

precision is very important when users need very accurate results.

Recall is a core effectiveness metrics used to measure how well an IR system retrieves all relevant documents.

It is the proportion of relevant documents that the system successfully retrieved.

$$\text{Recall} = \frac{\text{Relevant retrieved documents}}{\text{Total relevant documents.}}$$

Interpretation of all the relevant documents that exist how many did the system find.

High Recall. System retrieves most of the relevant documents.

Low Recall. System missed many relevant documents.

Example: Suppose in the entire db

- ✓ There are 20 relevant documents
- ✓ The system retrieves 12 relevant documents.

$$\text{Recall} = \frac{12}{20} = 0.6 = 60\%.$$

- Recall is more important when:
 - ✓ Missing a relevant document is dangerous or costly.

F1-score also called F1-Measure is a metric used to evaluate IR systems by combining precision and recall into a single score.

It is useful when both precision and recall are important and must be balanced.

It is the harmonic mean of precision and recall.

$$F1\text{-Score} = \frac{2 \times (\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}}$$

- # It gives higher score when both precision and recall are high.

Accuracy in IR System :- This is a metric that measures how often an IR system correctly classifies documents as relevant or non-relevant. It considers all four outcomes:-

$$\text{Accuracy}(A) = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

Example

Suppose an IR system processes 100 documents
 $\text{TP} = 30$ (Relevant retrieved) $\text{TN} = 60$ (Not relevant Not Retrieved)
 $\text{FP} = 5$ $\text{FN} = 5$

$$A = \frac{30 + 60}{30 + 60 + 5 + 5} = \frac{90}{100} = 0.9 = \underline{\underline{90\%}}$$

The Notion of Relevance in IR

- * Relevance is a central concept in information retrieval (IR) that measures how well a retrieved document meets the user's information need. A document is considered relevant if it satisfies the user's query intention, provides useful information, or helps solve the user's problem.
- Relevance is subjective, meaning different users may judge the same document differently depending on their knowledge, background, purpose and context.
- Relevance in IR is not just about matching keywords, it reflects how well the retrieved information aligns with the user's real info need, which may not always be fully expressed in the initial query. Because of this, IR systems often incorporate methods like Relevance feedback, Query expansion, and ranking algorithms to estimate and improve relevance.

Relevance can be understood in different forms, such as topical relevance (Content matches the topic), cognitive relevance (fits the user's understanding), situational relevance (supports the user's task), and contextual relevance (useful within a specific scenario). Measuring relevance is essential for evaluating IR system performance through metrics such as Precision, Recall, and F1-Score.

- ⇒ Relevance is a subjective judgment and may include:-
 - Being timely (recent information)
 - Being authoritative (from a trusted source)
 - Satisfying the goals of the user and his/her intended use of the information (information need).
 - Relevant information is that suited to your info need.
what is actually needed (relevant)
 - Dependent on: (user, space/time, Group and context)
- IR is very concerned with relevance.

Chapter 5

Retrieval Effectiveness

$$\text{Precision} = \frac{\text{Relevant Retrieved}}{\text{Total Retrieved}}$$

$$\text{Recall} = \frac{\text{Relevant Retrieved}}{\text{Total Relevant}}$$

$$\text{F1-score} = \frac{2(P \times R)}{P + R}$$

$$\text{Accuracy} = \frac{\text{Relevant Retrieved} + \text{Irrelevant NOT retrieved}}{\text{Total document repository}}$$

$$E\text{ measure} = 1 - \text{F1 score}$$

$$\text{Fallout} = \frac{\text{Retrieved Irrelevant}}{\text{Total Irrelevant}}$$

$$\text{Noise} = 1 - \text{Precision}$$

$$\text{Silence (Miss)} = 1 - \text{Recall}$$

Chapter 6:

Relevance feedback & Query Expansion

(Query operations)

* Query operations in IR are performed using different approaches to express user information needs more precisely and to improve retrieval effectiveness. The main approaches include:

- ✓ Boolean
- ✓ Vector Space
- ✓ Probabilistic.
- ✓ Fuzzy/ Approximate Methods

Boolean Approach

"Query operation is Actions on the query.
(Add, Remove, Reweight terms)

- ✓ Uses logical operators (AND, OR, NOT) to combine query terms.
- ✓ Document either matches or don't match the query
(Binary Relevance)
- ✓ Pros: Simple & precise for complex conditions
- ✓ Cons: No Ranking: Results are either retrieved or not, which may limit usefulness.

Example: Education and Technology NOT online

Vector Space Approach

- ✓ Represents queries and documents as vectors in a multidimensional space.
- ✓ Terms are assigned weights based on importance (e.g. TF-IDF)
- ✓ Documents are ranked based on similarity measures.
Usually Cosine Similarity.
- ✓ Pros: provides ranked results and partial matching
- ✓ Cons: Computationally more complex than Boolean.

Example: Query Vector: (0.2, 0.5, 0.3), document Vector (0.1, 0.4, 0.5)

* Relevance feedback (User, pseudo, interactive): Both a query operation (expanding/reforming the query) and an approach (manual, automatic, iterative) to improve retrieval.

⇒ Relevance feedback in IR is a special type of query operation where the system improves search results using information about which documents are relevant or not. It is a technique for query reformulation, expansion and an approach to improving retrieval effectiveness.

1. User Relevance feedback / Explicit feedback)

defn: The user explicitly marks retrieved documents as relevant or non-relevant.

Query operations involved:- Adds or reweights terms from relevant documents → query expansion/reformulation.

Approach: iterative refinement guided by actual user judgements.

Pros: Highly accurate because it reflects real user intent.

Cons: Requires user effort and time.

Example: After retrieving 20 articles, the user marks 5 as relevant; the system uses terms from these 5 to refine query.

2. Pseudo-Relevance feedback :- also called Blind feedback.

defn: The system assumes the top-ranked documents from the initial query are relevant without asking the user.

Query operations involved: Automatically Selects important terms from top documents → query expansion

Approach: System-driven, automatic refinement

Pros: No user effort required; can improve Recall

Cons: If top documents are irrelevant, precision may decrease

Example: Top 10 documents from a search on "Climate Change" are used to expand the query with terms like "Carbon emissions", "global warming" etc.

3. User's (interactive / iterative) Relevance feedback

* ^{defn:} Combine user input and iterative query refinement

Query operations involved: Query is reformulated based on user judgements multiple times → iterative expansion / term rewriting

Approach: interactive Method that repeatedly improves the query.

Pros: Highly effective; aligns with actual user intent

Cons: Requires multiple interactions and time.

Example: A user repeatedly marks relevant documents and the system adjusts the query after each iteration until results are satisfactory.

User Relevance feedback (URF)

Meaning: This is the traditional/classical relevance feedback process. The system retrieves a set of documents → the user marks which documents are relevant or not.

- The system updates the query.

Key Characteristics

Aspects

Who initiates?

User Role

System Role

interaction level

Focus

Description

- The system first presents results; the user provides feedback afterwards.
- User labels documents (relevant / not relevant)
- System adjusts the query using Rocchio / similar algorithms.
- Low to Moderate → user mainly gives labels.
- Improve precision and recall by modifying the query.

Example:
User Searches: "Water pollution Ethiopia"
System returns a list
User Marks: Relevant ✓ 5 documents
NOT Relevant ✗ 3 documents
System refines the query!
"Water Contamination, river pollution,
Environmental Management Ethiopia"

Interactive Relevance feedback (IRF)

Meaning: A more dynamic, iterative and user-system cooperative process. The system and the user work together through multiple rounds, not just one.

Key characteristics

Aspect

- ✓ Who initiates?
- ✓ User Role
- ✓ System's Role

Interaction level

Focus

Description

- ✓ Both user and system. User may directly modify query terms.
- User interacts with the interface: Select terms, expands query, browses clusters, adjust sliders, etc
- provides suggestions: term clouds, categories, filters, previouss.
- High - user controls the search process step by step
- Helping the user explore, refine, and understand the search space.

Core Difference

features

- ✓ interaction styles
- ✓ user engagement
- ✓ Query changes
- ✓ Rounds
- ✓ Goal
- ✓ Tools used

User

- ✓ one-way feedback
- ✓ Limited
- ✓ mostly automatic
- ✓ usually 1-2
- ✓ improves retrieval accuracy
- ✓ only document labeling

Interactive

- ✓ Two-way, continuous interaction
- ✓ High
- ✓ user + system together
- ✓ many rounds
- ✓ improve the search experience, exploration & control
- ✓ suggestion tools, filters, sliders, visualizations

In Simple terms :-

User Relevance feedback

"User marks documents as Relevant - system adjusts the query."

Mostly labeling

Interactive Relevance feedback

User and system work together step by step to refine the search.

Exploration + Collaboration + multiple interactions

⇒ User Relevance feedback

User: These 4 docs are Relevant

System: OK, I will modify the query

⇒ Interactive Relevance feedback

User: Add term: "River Contamination"

Remove term: "Pollution Control strategies"

Filter by year

Select cluster "Ethiopia Water Studies"

System Continuously updates

W

Techniques for enhancing Retrieval effectiveness in IR.

- * Improving retrieval effectiveness means improving how well an IR system retrieves relevant documents and filters out non-relevant ones.
 - The goal is to Maximize Precision, Recall and overall User satisfaction
- Below are the major techniques used to improve retrieval effectiveness

1. Query Reformulation is a core technique for enhancing Retrieval effectiveness.

Query Reformulation is the process of modifying, expanding or refining a user's original query to improve the quality of retrieval results.

- It helps ensure that the system retrieves documents that more accurately match the user's true information need.
- Why? Users often submit:

- ✓ Short queries ✓ ambiguous terms
- ✓ incomplete descriptions
- ✓ poorly expressed needs

Therefore, retrieval systems use reformulation to:-

- ✓ Clarify the query ✓ Add useful terms
- ✓ Remove noisy terms ✓ Improve precision & Recall

Techniques of Query Reformulation

= Techniques for enhancing retrieval effectiveness.

a) Manual Query Reformulation (User driven)

* The user modifies the query manually :-

Example include :-

- ✓ Using boolean operators [AND, OR & NOT]
- ✓ Adding Synonyms
- ✓ Removing irrelevant key words
- ✓ Narrowing & broadening Search terms
- Good when the user knows the domain

b) Automatic Query Reformulation (System driven)

* The system changes the query w/o user intervention

Common Methods :-

✓ Query expansion (QE)

Adding Synonyms, related terms, stems

Using thesauri (WordNet), ontologies, term co-occurrence

✓ Query Reduction

Removing Unhelpful or noisy terms

Keeping only the most discriminative terms.

✓ Spelling correction / Normalization

Correcting Spelling errors.

Handling Variants (Color → colour)

✓ Feedback-based Reformulation

User Relevance feedback (URF) Explicit

Pseudo-Relevance Feedback (Implicit)

Interactive Feedback (IRF)

Chapter 6 Relevance

Example of Query Reformulation

1. User Query "Virus"

System expands to :-

→ Virus + infection, Malware, antiviral, pathogen.

2. User Query : "Education technology Ethiopia"

Reformulated to :

Education technology + ICT in Education,
digital learning, Ethiopia, e-learning

* Query Formulation vs Query Reformulation
are two different stages in the search process
in IR.

+ Query Formulation is the process by which
a user creates the first / initial query to express
their information need in a search system.

Characteristics

- ✓ Done before any search results appear
- ✓ Expresses the user's understanding of the topic
- ✓ May be incomplete, vague or ambiguous
- ✓ often short (Average 2-4 words)

Activities involved

- ✓ Selecting keywords
- ✓ Choosing filters

- ✓ Using boolean operators
- ✓ Typing a natural-language question

* A Thesaurus in IR is a structured, controlled Vocabulary that lists Concepts and shows the relationship b/w them. (synonyms, broader terms, narrower terms, and Related terms)

- It helps both Users and IR Systems overcome problems of Vocabulary Mismatch, ambiguity, and inconsistent terminology.

A thesaurus in IR is a tool that provides :-

- ✓ Synonyms (eg. "Car" - "automobile")
- ✓ Broader Term (BT)
- ✓ Narrower Term (NT)
- ✓ Related term (RT)
- ✓ preferred descriptors or standard terms

* In IR, It serves as a bridge b/w the user language and the indexing language.

Why Thesauri improves Retrieved Effectiveness.

- ✓ Handling Synonymy
- ✓ Handling polysemy (multiple meaning)
- ✓ Standardized Vocabulary
- ✓ Supporting query expansion
- ✓ Guiding novice users

—