

Formula of BM 25 Score:-

score(D, Q) = \sum\_{i=1}^n IDF(q\_i) \cdot \frac{f(q\_i, D) \cdot (k\_1 + 1)}{f(q\_i, D) + k\_1 \cdot (1 - b + b \cdot \frac{|D|}{avgdl})}

IDF(q\_i) = \ln\left(\frac{N - n(q\_i) + 0.5}{n(q\_i) + 0.5} + 1\right)

= \ln\left(\frac{N - n(q\_i) + 0.5 + n(q\_i) + 0.5}{n(q\_i) + 0.5}\right)

= \ln\left(\frac{N + 1}{n(q\_i) + 0.5}\right)

Nomenclature:

- ★ score(D, Q) = Score of a doc | query
- ★ f(q\_i, D) = No of times q\_i occurs in doc D
- ★ |D| = Length of doc D [no of words in D]
- ★ avgdl = Avg length of docs

Formula of TF-IDF Score:-

tf(t, d) = \frac{f\_{t,d}}{\sum\_{t' \in d} f\_{t',d}}

idf(t, D) = \log \frac{N}{|\{d : d \in D \text{ and } t \in d\}|}

Score(D, Q) = \sum\_{t \in Q} tf(t, D) \times idf(t, D)

Nomenclature:

- ★ score(D, Q) = Score of a doc | query
- ★ f\_{t,D} = No of times t occurs in doc D
- ★ tf(t, D) = \frac{no of times t occurs in D}{Total no of words in D}

Intuitive explanation of the formula

Diagram showing three documents D with terms t. Document 1 has t, Document 2 has t, Document 3 has t. Term t is common to all.

① f\_{t,D} \uparrow \Rightarrow S\_{t,D} \uparrow

② Suppose t occurs in n\_t Docs, useless \frac{n\_t}{N} \uparrow, useful \frac{n\_t}{N} \downarrow \Rightarrow f\_{t,D} \downarrow

S\_{t,D} \uparrow, \frac{n\_t}{N} \downarrow \Rightarrow \log\left(\frac{N}{n\_t}\right) \uparrow

score(D, Q) = \sum\_{i=1}^n IDF(q\_i) \cdot \frac{f(q\_i, D) \cdot (k\_1 + 1)}{f(q\_i, D) + k\_1 \cdot (1 - b + b \cdot \frac{|D|}{avgdl})}

IDF(q\_i) = \ln\left(\frac{N - n(q\_i) + 0.5}{n(q\_i) + 0.5} + 1\right)

Problems with TF-IDF

① Diagram showing two documents A and B. Document A has terms t, Document B has terms t. Term t is common to both.

S\_c \approx 2 \times S\_D, S\_A \approx S\_B

Graph of f\_{t,D} vs tf(t, D). TF-IDF is shown as a curve that is higher than the expected curve.

② Doc Length bias: Longer Doc \rightarrow more chances of finding t \Rightarrow tf \uparrow \Rightarrow Score \uparrow \uparrow \uparrow

③ Lack of Flexibility in TFIDF: k\_1 & b \rightarrow Hyper parameters in BM25

k & b ARE FREE PARAMS. Typically: k \in [1.2, 2.0], b = 0.75

\frac{\partial BM25}{\partial \theta} = i \left[ \frac{-t(k+1)kb}{(t+k(1-b+b\theta))^2} \right]

BM25 \uparrow, \theta \downarrow

\* NEG REL. TO \theta

\* BUT DIMINISHING RETURNS TO \theta

\theta = \frac{|D|}{avgdl}

\frac{\partial BM25}{\partial t} = i \left[ \frac{k(k+1)(1-b+b\theta)}{(t+k(1-b+b\theta))^2} \right]

BM25 \uparrow, t \uparrow

TF-IDF

BM25

\* Pos REL. TO t

\* BUT DIMINISHING RETURNS TO t

f(q\_i, D)

BM25(2, 4) = 1

BM25(2, 8) = 0.2

BM25 (Best Matching 25) is a widely used ranking algorithm in information retrieval systems. It is employed in various applications and platforms:

Search Engines

- Web search engines like Google, Bing, and Yahoo use BM25 or similar algorithms to determine the relevance of search results3.
- Enterprise search systems in large organizations utilize BM25 to provide employees with relevant documents and information from internal databases3.

E-commerce

- Online shopping platforms often implement BM25 or similar algorithms to rank products based on relevance to user queries and provide personalized product recommendations3.

Information Retrieval Systems

- BM25 serves as a strong baseline in information retrieval research, particularly in the TREC Web track5.
- It is used in document retrieval systems to rank documents based on their relevance to search queries1.

Vector Databases

- Vector databases like Milvus integrate BM25 to enhance search relevance and efficiency1.

Hybrid Search Systems

- Many real-world search applications combine BM25-based search with vector-based semantic search powered by large language models (LLMs)4.

AI and Machine Learning

- BM25 is often integrated with LLMs using Retrieval-Augmented Generation (RAG) to improve search and retrieval performance4.
- It can be used as a cost-effective semantic cache when integrating LLMs into production systems4.

Specific Platforms

- Azure AI Search uses BM25 as its default relevance scoring algorithm11.
- Elasticsearch implements BM25 as its default similarity ranking algorithm10.
- SAP HANA Cloud Database includes BM25 search functionality9.
- LangChain and Weaviate, popular AI development frameworks, offer BM25 retrieval options1213.

By leveraging BM25 in these diverse applications, developers and researchers can create more effective and efficient information retrieval systems across various domains

Idf = \log\left(\frac{N}{N\_t}\right)

H = -\sum \log\left(\frac{1}{N}\right)

Surprise

\log\left(\frac{N\_t}{N}\right)