

Text Analysis and Retrieval

3. Basics of Information Retrieval

Prof. Jan Šnajder

University of Zagreb
Faculty of Electrical Engineering and Computing (FER)

Academic Year 2022/2023

1 Main IR models

2 IR evaluation

3 Neural IR

Reminder: What is IR?

Information retrieval

(Wikipedia)

The activity of obtaining **information resources** relevant to a user's **information need** from a collection of information resources.

- **Information needs** (expressed by users in the form of **queries**)
- **Information resources** (typically unstructured – text, images, video, audio, etc.)

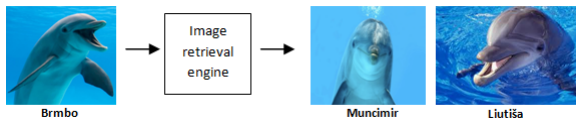
Information needs

Information need

Information need is an individual or group's **desire** to locate and obtain **information** to satisfy a conscious or unconscious **need**. Needs and interests call forth information.

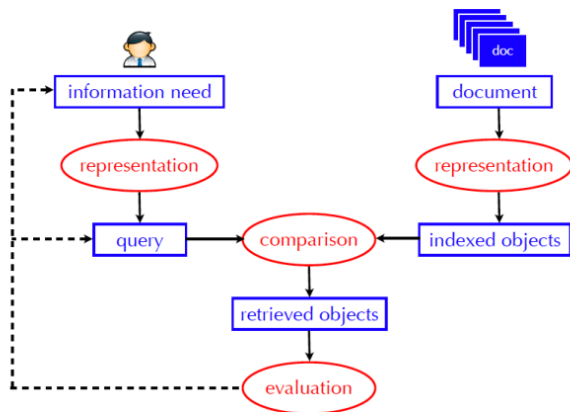
(Robert S. Taylor: "The process of asking questions", 2007)

- (Un)conscious needs for information are expressed via **queries**
 - In text retrieval: **words and phrases**
(e.g., "ISIS attacks", "coronavirus pandemic")
 - In image content retrieval: **images**



Information retrieval problem

Basic Information Retrieval Process



7

Diagram from Jaime Arguello, UNC-Chapel Hill

Unstructured document representation

Document snippet

One evening Frodo and Sam were walking together in the cool twilight. Both of them felt restless again. On Frodo suddenly the shadow of parting had falling: the time to leave Lothlorien was near.



Bag of words

{(One, 1), (evening, 1), (Frodo, 2), (and, 2), (Sam, 1) (were, 1), (walking, 1), (together, 1), (in, 1), (the, 3), (cool, 1), (twilight, 1), (Both, 1), (of, 2), (them, 1), (felt, 1), (restless, 1), (again, 1), (On, 1), (suddenly, 1), (shadow, 1), (parting, 1), (had, 1), (falling, 1), (time, 1), (to, 1), (leave, 1), (Lothlorien, 1), (was, 1), (near, 1)}

Weakly-structured document representation

Document snippet

One evening Frodo and Sam were walking together in the cool twilight. Both of them felt restless again. On Frodo suddenly the shadow of parting had falling: the time to leave Lothlorien was near.



Bag of nouns

{(evening, 1), (Frodo, 2), (Sam, 1), (twilight, 1), (shadow, 1), (parting, 1), (time, 1), (Lothlorien, 1)}

Bag of named entity terms

{(Frodo, 2), (Sam, 1), (Lothlorien, 1)}

① Morphological normalization (stemming/lemmatization)

- Conflating various forms of the same word to a common form
- Important for morphologically rich languages such as Croatian
- **Stemming** (e.g., kućom → kuć) more often used than **lemmatization** (e.g., kućom → kuća)

② Stop words removal

- **Stop words**: semantically void terms such as determiners, prepositions, conjunctions, pronouns, etc.
- **Content words**: nouns, verbs, adjectives, adverbs
- Stop words removal: removes stop words and retains only the content words
- English stop words: <https://www.ranks.nl/stopwords>

⇒ both methods reduce the size of the bag-of-words representation and generally improve retrieval performance

Three components of an IR model

The **basic retrieval model** is a triple (f_d, f_q, r) :

- 1 f_d is a function that maps **documents** to retrieval representations

$$f_d(d) = x_d$$

- 2 f_q is a function that maps **queries** to retrieval representations

$$f_q(q) = x_q$$

- 3 r is a **ranking function** that produces a **relevance score**: a real-valued number that indicates the potential relevance of the document d for query q based on x_d and x_q

$$relevance(d, q) = r(f_d(d), f_q(q)) = r(x_d, x_q)$$

- Information retrieval models roughly fall into three paradigms:
 - 1 **Set-theoretic models**
 - Boolean model
 - Extended Boolean model
 - 2 **Algebraic models**
 - Vector space model
 - Latent semantic indexing
 - 3 **Probabilistic models**
 - BM25
 - Language model
- Additionally, there are IR models that utilize **link analysis algorithms** (e.g., PageRank, HITS)

Vector space model

- Documents and queries are represented as vectors of index terms
- Weights are non-negative real numbers

$$\mathbf{d}_j = [w_{1j}, w_{2j}, \dots, w_{tj}]$$

$$\mathbf{q} = [w_{1q}, w_{2q}, \dots, w_{tq}]$$

- The relevance of the document for the query is estimated by computing some **distance or similarity metric** between the two vectors
 - Distance metrics – Euclidean, Manhattan, etc.
 - More relevant when distance is lower
 - Similarity metrics – Cosine, Dice, etc.
 - More relevant when similarity is larger

Cosine similarity

Salton, 1983

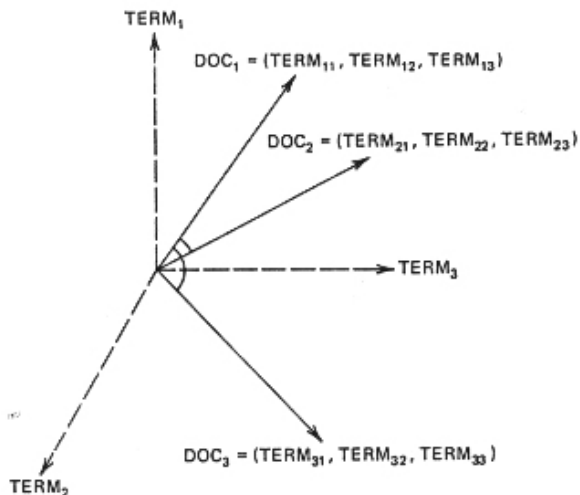


Figure 4-2 Vector representation of document space.

Vector space model – term weighting

- How are the weights w_{ij} of index terms for documents computed?
- Two intuitive assumptions:
 - ① The relevance of an index term for the document is proportional to its frequency in the document (**term frequency** component)
 - i.e., **more frequent** \Rightarrow **more relevant**
 - ② The relevance of an index term for any document is inversely proportional to the number of documents in the collection in which it occurs (**inverse document frequency** component)
 - i.e., **more common across documents** \Rightarrow **less relevant**
(e.g., stopwords such as “*the*”)

TF-IDF weighting scheme

- The weight computed as the product of the term frequency component and the inverse document frequency component

$$w_{ij} = tf(t_i, d_j) \cdot idf(t_i, D)$$

where t_i is the i -th term from the index

- The most popular local and global schemes:

$$tf(t_i, d_j) = 0.5 + \frac{0.5 \cdot freq(t_i, d_j)}{\max_{t \in d_j} freq(t, d_j)}$$

$$idf(t_i, D) = \log \frac{|D|}{|\{d_j \in D \mid t_i \in d_j\}|}$$

- **Probabilistic retrieval models**
 - View retrieval as a problem of **estimating the probability of relevance** given a query, document, collection, etc.
 - Documents are ranked in decreasing order of this probability
- Rely on **probability ranking principle**

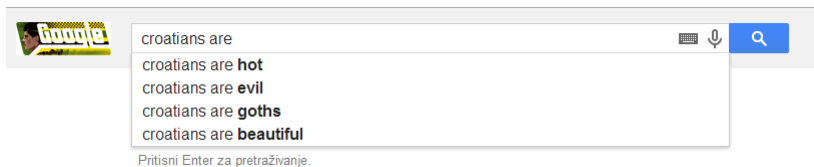
- The famous **BM25 ranking function**:

$$\sum_{t \in q} \frac{\text{freq}(t, d)(k_1 + 1)}{k_1(1 - b) + k_1(l_d/l_{avg})b + \text{freq}(t, d)} \cdot \text{idf}(t, d)$$

- Often yields SOTA results!

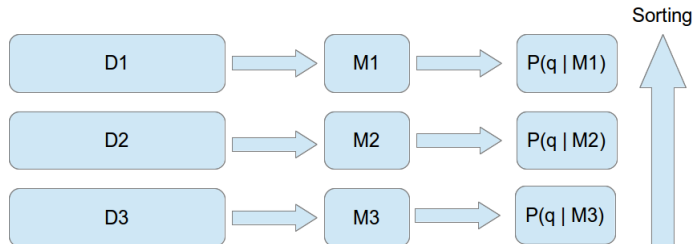
Language modeling for IR

- Approaching the probabilistic information retrieval problem from a different perspective
- Instead of modeling document probability given the query, we model the **query probability given the document**



Query likelihood model

- Given a document collection D and a query q
- A language model M_d is built for **each document**
- Documents are scored according to the probability $P(q|M_d)$



- **Intuition:** language models corresponding to relevant documents should assign higher probability to the query

- **Terrier IR platform**

- Open command-line IR platform
- Popular in academia – used for IR evaluations in research
- Stepwise usage – first indexing, then retrieval, and then evaluation

- **Lucene**

- Open source information retrieval library written in Java (ports to many languages exist)
- Describes a document as a set of (user definable) fields
- Very flexible solution for applications requiring full text indexing and search

- **Elasticsearch** – RESTful search engine on top of Lucene

- Docs: <https://www.elastic.co/guide/index.html>
- Very limited NLP functionality:
<https://www.elastic.co/blog/text-classification-made-easy-with-elasticsearch>

Learning outcomes 1

- ① List three main components of an IR model
- ② Describe the vector space model and the TF-IDF weighting scheme
- ③ Explain the probability ranking principle and BM25
- ④ Describe the LM information retrieval model
- ⑤ List the main IR tools available

1 Main IR models

2 IR evaluation

3 Neural IR

IR evaluation

Clash of the giants: Which one is better?

The image shows a side-by-side comparison of search results for the query "kyoto public transportation". On the left is the Bing interface, and on the right is the Google interface.

Bing Results (Left):

- Search bar: "bing vs. Google kyoto public transportation" with buttons for "Search", "horizontal split", "bing only", "google only", and "add to browser".
- Navigation: WEB, IMAGES, VIDEOS, NEWS, MORE.
- Results: 656,000 RESULTS.
- Top results include:
 - Kyoto City Web / Access / Public transport in Kyoto** (www.city.kyoto.jp/koho/eng/access/transport.html) - The Kyoto City bus is useful for getting around various places within Kyoto. Most of the city buses look like the diagram below: About the City Bus
 - Kyoto: Public Transportation - TripAdvisor** (www.tripadvisor.com/.../Kyoto:Japan:Public.Transportation.html) - 3/20/2014 - Inside Kyoto: Public Transportation - Before you visit Kyoto, visit TripAdvisor for the latest info and advice, written for travelers by travelers.
 - Kyoto City Web / Access** (www.city.kyoto.jp/koho/eng/access/index.html) - Public transport in Kyoto : Subway Map : Topics of City Government : Tourist Info : Useful Living Info : Kyoto Game Watching : Site Map : Link Info: Access: Access to ...
 - Transportation in Japan** (www.japan-guide.com/e/e627.html) - Japan has an efficient public transportation network, especially within metropolitan areas and between the large cities. Japanese public transportation is ...
 - Kyoto Travel: Access, Orientation and Transportation** (www.japan-guide.com/e/e2363.html) - Kyoto has a rather inadequately developed public transportation system for ... Itocha and Pitapa can be used on most means of public transportation in Kyoto and ...
 - Kyoto Public Transport guide and map. - HotelTravel.com** (www.hoteltravel.com/tourism/kyoto/kyoto-public-transport.html) -

Google Results (Right):

- Search bar: "Search Images Maps Play YouTube News Gmail Drive More" with "Sign In" button.
- Navigation: Web, Images, Videos, News, Shopping, Maps, Books.
- Results: About 1,070,000 results.
- Top results include:
 - Any time** (Past hour, Past 24 hours, Past week, Past month, Past year)
 - Kyoto City Web / Access / Public transport in Kyoto** (https://www.city.kyoto.jp/koho/eng/access/transport.html) - The Kyoto City bus is useful for getting around various places within Kyoto. Most of the ... Please enter the bus from the back door, and exit at the front. The bus ...
Subway Map - Access to Kyoto City - Tourist Info
 - Kyoto Travel: Access, Orientation and Transportation - Japan Guide** (www.japan-guide.com/e/e2363.html) - The closest airport to Kyoto is Osaka's Itami Airport, about one hour by bus from central Kyoto (more details). Most flights connect Itami Airport with Tokyo's ...
 - Kyoto Visitor's Guide-Transportation System-** (www.kyotoguide.com/ver2/useful/useful-trans.htm) - In Kyoto, you enter the bus from the back, exit and pay at the front. Change for 500 yen and 1,000 yen bills, etc. can be made by the machine at the front of ...
 - Useful Services & Tickets - Kyoto Travel Guide** (www.kyoto.travel/2009/11/useful-services-tickets.html) -
C. City Bus and Subway Information Center (Chikatsutsu Annai-sho) in Kotochika Kyoto (Delivered by Yamato Unyu) Tel: +81-(0)75-371-9866. Open 7:30 to 12: ...

IR test collection

- IR test collection is comprised of:
 - ① Document collection
 - ② Set of information needs (descriptions + queries)
 - At least 50 information needs
 - ③ Set of relevance judgments for each query–document pair
 - Binary relevance (document is **relevant** or **not relevant**) or graded relevance judgments (less common)
- Used for:
 - Evaluating retrieval effectiveness w.r.t. different settings (stemming, ranking model, etc.)
 - Comparing against other systems, typically in **evaluation campaigns**
 - Fine-tuning of system parameters, done on a **development test collection** (to prevent overfitting)

IR test collections – topic examples

TREC topic 351

Title: Falkland petroleum exploration

Description: What information is available on petroleum exploration in the South Atlantic near the Falkland Islands?

Narrative: Any document discussing petroleum exploration in the South Atlantic near the Falkland Islands is considered relevant. Documents discussing petroleum exploration in continental South America are not relevant.

TREC topic 409

Title: Legal, Pan Am, 103

Description: What legal sanctions have resulted from the destruction of Pan Am Flight 103 over Lockerbie, Scotland, on December 21, 1988?

Narrative: Documents describing any charges, claims, or fines presented to or imposed by any court or tribunal are relevant, but documents that discuss charges made in diplomatic jousting are not relevant.

Relevance judgments (“qrel file”)

401 0 FBIS3-18916 0
401 0 FBIS3-18926 0
401 0 FBIS3-18943 1
401 0 FBIS3-18946 0
401 0 FBIS3-18972 0
401 0 FBIS3-18997 0
401 0 FBIS3-19003 0
401 0 FBIS3-19032 1
401 0 FBIS3-19037 0
401 0 FBIS3-19038 1
401 0 FBIS3-19042 0
401 0 FBIS3-19080 0
401 0 FBIS3-19103 0
401 0 FBIS3-19107 1
401 0 FBIS3-19110 0
401 0 FBIS3-19126 0
401 0 FBIS3-19133 0
401 0 FBIS3-19212 0
401 0 FBIS3-19213 0
401 0 FBIS3-19251 0
401 0 FBIS3-19290 0
401 0 FBIS3-19302 0
401 0 FBIS3-19303 0
401 0 FBIS3-19304 0

Evaluation metrics

- Compare retrieved documents against relevant documents
- Each document is either retrieved or not, and either relevant or not. This gives a 2×2 confusion matrix:

	relevant	not relevant
retrieved	tp	fp
not retrieved	fn	tn

We could compute **accuracy** as the fraction of correct decisions:

$$Acc = \frac{tp + tn}{tp + tn + fp + fn}$$

- **Q:** Why using accuracy is not a good idea?
A: Given a query, most documents (say 99%) are irrelevant. A search engine that retrieves nothing will already be 99% accurate

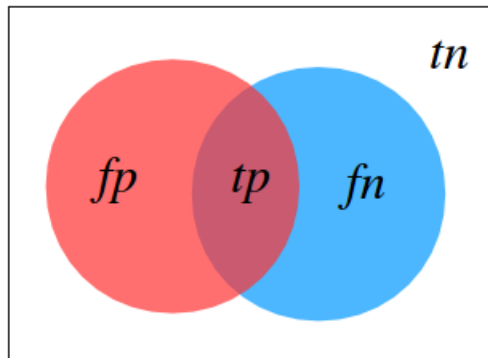
- **Precision (P)** is a fraction of retrieved documents that are relevant

$$P = \frac{\#(\text{relevant documents retrieved})}{\#(\text{retrieved documents})} = \frac{tp}{tp + fp}$$

- **Recall (R)** is a fraction of relevant documents that are retrieved

$$R = \frac{\#(\text{relevant documents retrieved})}{\#(\text{relevant documents})} = \frac{tp}{tp + fn}$$

Precision and Recall



$$P = \frac{tp}{tp + fp} \quad R = \frac{tp}{tp + fn}$$

F-measure

- Combining P and R into a single number (the harmonic mean)

$$F = \frac{1}{\alpha \frac{1}{P} + (1 - \alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R} \quad \text{where } \beta^2 = \frac{1 - \alpha}{\alpha}$$

- β controls the **precision–recall trade-off**:
 - $\beta = 1$ gives equal weight to precision and recall (**F1-score**):

$$F_{\beta=1} = \frac{2PR}{P + R}$$

- $\beta = 0.5$ emphasizes precision twice as much as recall
- $\beta = 2$ emphasized recall twice as much as precision

Evaluation of ranked results

- Modern search engines produce **ranked results**, but P , R , and F -score do not account for ranks
- Ideally, a search engine should rank all the relevant documents before the non-relevant ones
- Rank-based metrics:
 - Precision-recall curve
 - 11-point precision
 - MAP
 - $P@k$
 - R-precision
 - MRR

Learning outcomes 2

- 1 Explain what an IR test collection consist of and what it's used for
- 2 Define and calculate the standard IR evaluation metrics

1 Main IR models

2 IR evaluation

3 Neural IR

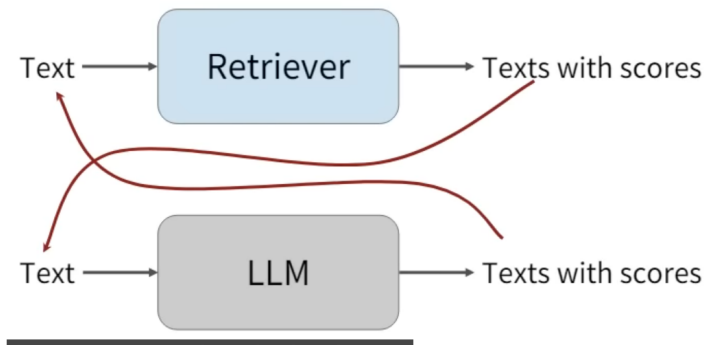
- BERT+search:

- <https://jalammr.github.io/illustrated-retrieval-transformer/>
- <https://blog.google/products/search/search-language-understanding-bert/>
- [https://azure.microsoft.com/en-us/blog/bing-delivers-its-largest-improvement-in-search-experience-using-azure-g](https://azure.microsoft.com/en-us/blog/bing-delivers-its-largest-improvement-in-search-experience-using-azure-gpt-4/)

- Generative LLMs for IR (also: conversational search):

- <https://www.technologyreview.com/2021/05/14/1024918/language-models-gpt3-search-engine-google/>
- Metzler et al. (2021). Rethinking search: making domain experts out of dilettantes. ACM SIGIR.
<https://arxiv.org/pdf/2105.02274.pdf>
- Shah & Bender. (2022). Situating search. ACM SIGIR.
<https://dl.acm.org/doi/abs/10.1145/3498366.3505816>

Models can communicate in natural language



Christopher Potts: Stanford Webinar – GPT-3 & Beyond.

<https://www.youtube.com/watch?app=desktop&v=-lnHHWRCDGk>

Learning outcomes 3

- 1 Explain LLM-based search and what the advantages of this approach are over standard IR
- 2 Summarize the main critical points against LLM-based search and exemplify search scenarios in which such search fails

Study assignment

- 1 Watch TAR “Basics of IR 1/3” video:

<https://www.youtube.com/watch?v=MtzeLh99VHc&feature=youtu.be>

- 2 Watch TAR “Basic of IR 2/3” video:

<https://www.youtube.com/watch?v=DeMc-0oGem8&feature=youtu.be>

- 3 Read these three articles on using LLMs for search:

- <https://www.technologyreview.com/2021/05/14/1024918/language-models-gpt3-search-engine-google/>
- <https://www.technologyreview.com/2023/02/16/1068695/chatgpt-chatbot-battle-search-microsoft-bing-google>
- <https://blogs.bing.com/search/february-2023/The-new-Bing-Edge-%E2%80%93-93-Learning-from-our-first-week>

- 4 Read (Shah & Bender, 2022) (focus on sections 1, 2, 4.2, and 4.3):

<https://dl.acm.org/doi/abs/10.1145/3498366.3505816>

- 5 Self-check against learning outcomes!