Chapter IR:V

V. Retrieval Models

- Overview of Retrieval Models
- Empirical Models
- Boolean Retrieval
- Vector Space Model
- Probabilistic Models
- □ Binary Independence Model
- □ Okapi BM25
- Hidden Variable Models
- Latent Semantic Indexing
- □ Explicit Semantic Analysis
- Generative Models
- Language Models
- □ Combining Evidence
- Web Search
- □ Learning to Rank

Empirical Models [Probabilistic Models] [Hidden Variable Models] [Generative Models]

Basic empirical retrieval models abstract over a document $d \in D$ by treating them as a "bag of words" comprising the index terms derived from d.

A document representation d is composed of weighted index terms of d.

Discriminating factors of empirical models:

- 1. Term weighting method to compute the weight w_i of an index term t_i .
- 2. Construction method of the query representation q.
- 3. Computation method of the relevance function $\rho(\mathbf{q}, \mathbf{d})$.
- 4. Composition method of the result set *R*.

Retrieval Model $\mathcal{R} = \langle \mathbf{D}, \mathbf{Q}, \rho \rangle$ [Generic Model] [VSM] [BIM] [BM25] [LSI] [ESA] [LM]

Document representations D.

The set of index terms $T = \{t_1, \dots, t_m\}$ is composed of nouns as lemmatized word stems.

The representation d of a document d is a function from T to $\{0,1\}$, where $\mathbf{d}(t_i) = 1$ is interpreted as "term t_i present in d", and $\mathbf{d}(t_i) = 0$ as "term t_i absent from d".

Query representations Q.

A query representation \mathbf{q} corresponds to a logical formula with alphabet $\Sigma = T$, where the logical operators \wedge , \vee , \neg , and brackets can be used.

Relevance function ρ .

The document representation \mathbf{d} of a document d induces an interpretation $\mathcal{I}_{\mathbf{d}}$ for \mathbf{q} , yielding $\rho(\mathbf{q}, \mathbf{d}) = \mathcal{I}_{\mathbf{d}}(\mathbf{q})$.

If $\rho(\mathbf{q}, \mathbf{d}) = 1$, the document d is an element of the result set R.

Retrieval Model $\mathcal{R} = \langle \mathbf{D}, \mathbf{Q}, \rho \rangle$ [Generic Model] [VSM] [BIM] [BM25] [LSI] [ESA] [LM]

Document representations D.

The set of index terms $T = \{t_1, \dots, t_m\}$ is composed of nouns as lemmatized word stems.

The representation d of a document d is a function from T to $\{0,1\}$, where $d(t_i) = 1$ is interpreted as "term t_i present in d", and $d(t_i) = 0$ as "term t_i absent from d".

Query representations Q.

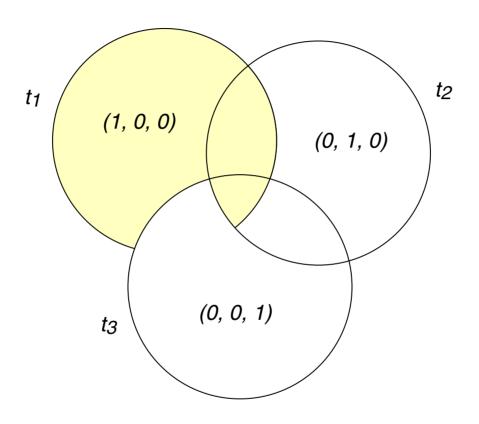
A query representation q corresponds to a logical formula with alphabet $\Sigma = T$, where the logical operators \land , \lor , \neg , and brackets can be used.

Relevance function ρ .

The document representation \mathbf{d} of a document d induces an interpretation $\mathcal{I}_{\mathbf{d}}$ for \mathbf{q} , yielding $\rho(\mathbf{q},\mathbf{d})=\mathcal{I}_{\mathbf{d}}(\mathbf{q}).$

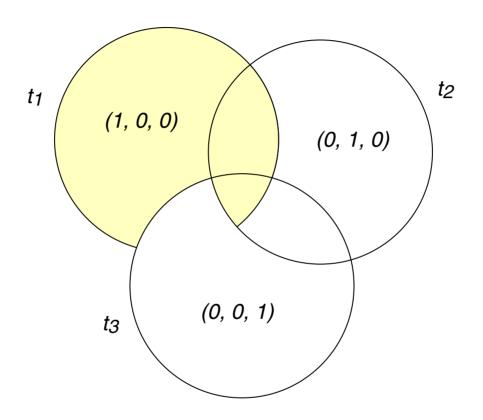
If $\rho(\mathbf{q}, \mathbf{d}) = 1$, the document d is an element of the result set R.

Relevance Function ρ



What is the query illustrated?

Relevance Function ρ



What is the query illustrated?

$$\mathbf{q} = t_1 \wedge (t_2 \vee \neg t_3) \approx (t_1 \wedge \neg t_2 \wedge \neg t_3) \vee (t_1 \wedge t_2 \wedge \neg t_3) \vee (t_1 \wedge t_2 \wedge t_3)$$

Example

Document representation:

$$\mathbf{d} = \{ \text{ (chrysler, 1), (deal, 1),} \\ \text{(usa, 1), (china, 0),} \\ \text{(cat, 0), (sales, 1),} \\ \text{(dog, 0), ...} \}$$

Query representation:

$$\mathbf{q} = \mathbf{usa} \wedge (\mathbf{dog} \vee \neg \mathbf{cat})$$
 $\approx (\mathbf{usa} \wedge \mathbf{dog}) \vee (\mathbf{usa} \wedge \neg \mathbf{cat})$
 $\approx (\mathbf{usa} \wedge \neg \mathbf{dog} \wedge \neg \mathbf{cat}) \vee$
 $(\mathbf{usa} \wedge \mathbf{dog} \wedge \neg \mathbf{cat}) \vee$
 $(\mathbf{usa} \wedge \mathbf{dog} \wedge \neg \mathbf{cat})$

Induces interpretation:

$$\mathcal{I}_{\mathbf{d}}(\mathbf{q}) = 1$$
, since $\mathcal{I}_{\mathbf{d}}(\mathtt{usa}) = 1$, $\mathcal{I}_{\mathbf{d}}(\mathtt{dog}) = 0$, and $\mathcal{I}_{\mathbf{d}}(\mathtt{cat}) = 0$.

Remarks:

- \Box The symbol " \approx " denotes "is logically equivalent with".
- What does logical equivalence mean?
- □ A Boolean query in disjunctive normal form can be answered straightforward using an inverted index in parallel for each conjunction.
- □ A Boolean query in canonical disjunctive normal will retrieve each document only once.

Query Refinement: "Searching by Numbers"

Best practice in Boolean retrieval: (re)formulate queries until the number of documents retrieved is manageable. Example: pages about President Lincoln.

- 1. lincoln
 Result: many pages about cars, places, people
- 2. president \(\) lincoln
 Result: "Ford Motor Company today announced that Darryl Hazel
 will succeed Brian Kelley as president of Lincoln Mercury."
- 3. president ∧ lincoln ∧ ¬automobile ∧ ¬car

 Not in result: "President Lincoln's body departs Washington in a nine-car funeral train."
- **4.** president \land lincoln \land ¬automobile \land biography \land life \land birthplace \land gettysburg Result: \emptyset
- 5. president \(\) lincoln \(\) \(\) automobile \(\) (biography \(\) life \(\) birthplace \(\) gettysburg)
 Top result might be: "President's Day Holiday activities crafts, mazes, word searches, ... 'The Life of Washington' Read the entire book online! Abraham Lincoln Research Site"

Query Refinement: "Searching by Numbers"

Best practice in Boolean retrieval: (re)formulate queries until the number of documents retrieved is manageable. Example: pages about President Lincoln.

- lincoln
 Result: many pages about cars, places, people
- 2. president ∧ lincoln
 Result: "Ford Motor Company today announced that Darryl Hazel
 will succeed Bria Kelley as president of <u>Lincoln</u> Mercury."
- 4. president \wedge automobile \wedge biography \wedge life \wedge

WAAAAHHH

ident (biography V life V

Day – Holiday activities – crafts, mazes, word searches, ad the entire book online! Abraham Lincoln Research Site"

Discussion

Advantages:

- Precision: in principle, any subset of documents from a collection can be designated by a Boolean query
- as in data retrieval, other fields are possible (e.g, date, document type, etc.)
- simple, efficient implementation

Disadvantages:

- retrieval effectiveness depends entirely on the user
- cumbersome query formulation (e.g., expertise required)
- no possibility to weight query terms
- no ranking; binary relevance soring is too restrictive for practical purposes
- the size of the result set is difficult to be controlled

Retrieval Model $\mathcal{R} = \langle \mathbf{D}, \mathbf{Q}, \rho \rangle$ [Generic Model] [Boolean Retrieval] [BIM] [BM25] [LSI] [ESA] [LM]

Document representations D.

The set of index terms $T = \{t_1, \dots, t_m\}$ is typically composed of the word stems of the vocabulary of a document collection, excluding stop words.

The representation d of a document d is a |T|-dimensional vector, where the i-th vector component of d corresponds to a term weight w_i of term $t_i \in T$, indicating its importance for d. Various term weighting schemes have been proposed.

Query representations Q.

A query representation q is constructed like a document representation

Relevance function ρ .

Document representations and query representations are interpreted as points in a vector space spanned by unit vectors for each term in T, assuming their orthogonality.

Distance and similarity functions defined for vector spaces serve as relevance functions ρ . The Euclidean distance and the cosine similarity are important examples.

Retrieval Model $\mathcal{R} = \langle \mathbf{D}, \mathbf{Q}, \rho \rangle$ [Generic Model] [Boolean Retrieval] [BIM] [BM25] [LSI] [ESA] [LM]

Document representations D.

The set of index terms $T = \{t_1, \dots, t_m\}$ is typically composed of the word stems of the vocabulary of a document collection, excluding stop words.

The representation d of a document d is a |T|-dimensional vector, where the i-th vector component of d corresponds to a term weight w_i of term $t_i \in T$, indicating its importance for d. Various term weighting schemes have been proposed.

Query representations Q.

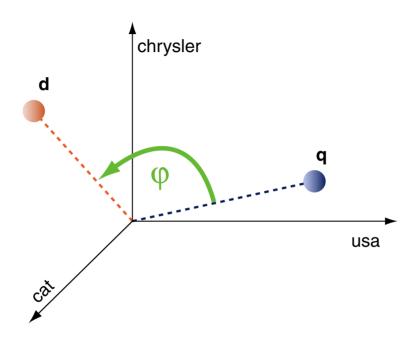
A query representation q is constructed like a document representation.

Relevance function ρ .

Document representations and query representations are interpreted as points in a vector space spanned by unit vectors for each term in T, assuming their orthogonality.

Distance and similarity functions defined for vector spaces serve as relevance functions ρ . The Euclidean distance and the cosine similarity are important examples.

Relevance Function ρ : Cosine Similarity



Relevance Function ρ : Cosine Similarity

The scalar product $\mathbf{a}^T\mathbf{b}$ between two vectors \mathbf{a} and \mathbf{b} , where φ denotes the angle between them, is defined as follows:

$$\mathbf{a}^{T}\mathbf{b} = ||\mathbf{a}|| \cdot ||\mathbf{b}|| \cdot \cos(\varphi)$$

$$\Leftrightarrow \cos(\varphi) = \frac{\mathbf{a}^{T}\mathbf{b}}{||\mathbf{a}|| \cdot ||\mathbf{b}||}$$

Relevance Function ρ : Cosine Similarity

The scalar product $\mathbf{a}^T\mathbf{b}$ between two vectors \mathbf{a} and \mathbf{b} , where φ denotes the angle between them, is defined as follows:

$$\mathbf{a}^{T}\mathbf{b} = ||\mathbf{a}|| \cdot ||\mathbf{b}|| \cdot \cos(\varphi)$$

$$\Leftrightarrow \cos(\varphi) = \frac{\mathbf{a}^{T}\mathbf{b}}{||\mathbf{a}|| \cdot ||\mathbf{b}||}$$

Normalizing a and b, denoted as a' and b', yields:

$$\cos(\varphi) = \frac{\mathbf{a}^T \mathbf{b}}{||\mathbf{a}|| \cdot ||\mathbf{b}||} = \frac{(\mathbf{a}')^T \mathbf{b}'}{||\mathbf{a}'|| \cdot ||\mathbf{b}'||} = (\mathbf{a}')^T \mathbf{b}' = \sum_{i=1}^m a_i' \cdot b_i'$$

 $\mathcal{R} = \langle \mathbf{D}, \mathbf{Q}, \rho \rangle$ with cosine similarity:

Definition of $\rho(\mathbf{q}, \mathbf{d})$ as $\cos(\varphi)$, where φ denotes the angle between \mathbf{q} and \mathbf{d} .

Example

$$\mathbf{d} = \begin{pmatrix} \mathsf{chrysler} & w_1 \\ \mathsf{usa} & w_2 \\ \mathsf{cat} & w_3 \\ \mathsf{dog} & w_4 \\ \mathsf{mouse} & w_5 \end{pmatrix} \quad = \quad \begin{pmatrix} \mathsf{chrysler} & 1 \\ \mathsf{usa} & 4 \\ \mathsf{cat} & 3 \\ \mathsf{dog} & 7 \\ \mathsf{mouse} & 5 \end{pmatrix}$$

$$\mathbf{d}' = \begin{pmatrix} \text{chrysler } 0.1 \\ \text{usa} & 0.4 \\ \text{cat} & 0.3 \\ \text{dog} & 0.7 \\ \text{mouse} & 0.5 \end{pmatrix} \text{,} \qquad \mathbf{q}' = \begin{pmatrix} \text{chrysler } 0.5 \\ \text{usa} & 0.5 \\ \text{cat} & 0.5 \\ \text{dog} & 0.5 \\ \text{elephant } 0.5 \end{pmatrix}$$

The angle φ between \mathbf{d}' and \mathbf{q}' is about 41° , $\cos(\varphi) \approx 0.75$.

Term Weighting: $tf \cdot idf$ [BIM Relevance Function]

To compute the weight w for a term t from document d under the vector space model, the most commonly employed term weighting scheme $\omega(t)$ is $tf \cdot idf$:

- 1. In document d, the importance of term t is proportional to its frequency.
 - \rightarrow term frequency tf(t, d) denotes the frequency of term t in document d. The term frequency is typically normalized, e.g., by the length of d.
- 2. In a document collection D, the importance of term t is inversely proportional to its document frequency df(t, D), the number of documents that contain t.
 - → inverse document frequency idf(t, D):

$$idf(t,D) = \log \frac{|D|}{df(t,D)}$$

The base of the logarithm is less important.

A term weight w for term t in document $d \in D$ is computed as follows

$$\omega(t) = tf(t, d) \cdot idf(t, D).$$

Term weights are normalized by $||\mathbf{d}||$ before indexing to simplify query processing.

Term Weighting: *tf* · *idf* [BIM Relevance Function]

To compute the weight w for a term t from document d under the vector space model, the most commonly employed term weighting scheme $\omega(t)$ is $tf \cdot idf$:

- 1. In document d, the importance of term t is proportional to its frequency.
 - \rightarrow term frequency tf(t, d) denotes the frequency of term t in document d. The term frequency is typically normalized, e.g., by the length of d.
- 2. In a document collection D, the importance of term t is inversely proportional to its *document frequency* df(t, D), the number of documents that contain t.
 - \rightarrow inverse document frequency idf(t, D):

$$idf(t,D) = \log \frac{|D|}{df(t,D)}$$

The base of the logarithm is less important.

A term weight w for term t in document $d \in D$ is computed as follows:

$$\omega(t) = tf(t, d) \cdot idf(t, D).$$

Term weights are normalized by $||\mathbf{d}||$ before indexing to simplify query processing.

Term Weighting: *tf* · *idf* [BIM Relevance Function]

To compute the weight w for a term t from document d under the vector space model, the most commonly employed term weighting scheme $\omega(t)$ is $tf \cdot idf$:

- 1. In document d, the importance of term t is proportional to its frequency.
 - \rightarrow term frequency tf(t, d) denotes the frequency of term t in document d. The term frequency is typically normalized, e.g., by the length of d.
- 2. In a document collection D, the importance of term t is inversely proportional to its *document frequency* df(t, D), the number of documents that contain t.
 - \rightarrow inverse document frequency idf(t, D):

$$idf(t,D) = \log \frac{|D|}{df(t,D)}$$

The base of the logarithm is less important.

A term weight w for term t in document $d \in D$ is computed as follows:

$$\omega(t) = tf(t, d) \cdot idf(t, D).$$

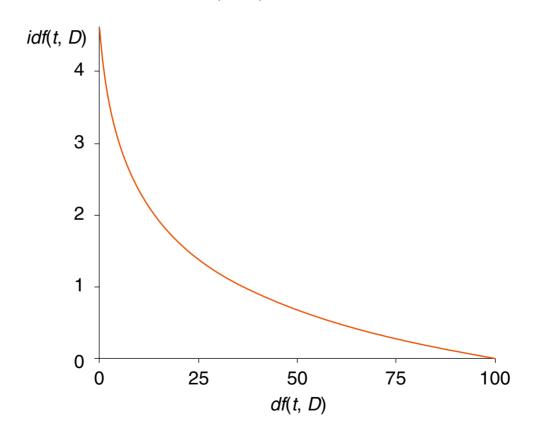
Term weights are normalized by $||\mathbf{d}||$ before indexing to simplify query processing.

Remarks:

- Term frequency weighting was invented by Hans Peter Luhn: "There is also the probability that the more frequently a notion and combination of notions occur, the more importance the author attaches to them as reflecting the essence of his overall idea." [Luhn 1957]
- \Box The importance of a term t for a document d is not linearly correlated with its frequency. Several normalization factors have been proposed [Wikipedia]:
 - tf(t,d)/|d|
 - $1 + \log(tf(t,d))$ for tf(t,d) > 0
 - $k + (1-k)\frac{tf(t,d)}{\max_{t' \in d}(tf(t',d))}$, where k serves as smoothing term; typically k = 0.4
- Inverse document frequency weighting was invented by Karen Spärck Jones: "it seems we should treat matches on non-frequent terms as more valuable than ones on frequent terms, without disregarding the latter altogether. The natural solution is to correlate a term's matching value with its collection frequency." [Spärck Jones 1972]
- □ Spärck Jones gives little theoretical justification for her intuition. Given the success of *idf* in practice, over the decades, numerous attempts at a theoretical justification have been made. A comprehensive overview has been compiled by [Robertson 2004].
- For example, interpreting the term $\frac{|D|}{df(t,D)}$ as inverse of the probability $P_{df}(t) = \frac{df(t,D)}{|D|}$ of t occurring in a random document in D yields $idf(t,D) = \log \frac{|D|}{df(t,D)} = -\log P_{df}(t)$. Logarithms fit relevance functions ρ since both are additive, yielding the interpretation: "The less likely (on a random basis) it is that a given combination of terms occurs, the more likely it is that a document containing this combination is relevant to the question." [Robertson 1972]

Term Weighting: tf · idf

Plot of the function
$$\mathit{idf}(t,D) = \log \frac{|D|}{\mathit{df}(t,D)}$$
 for $|D| = 100$.

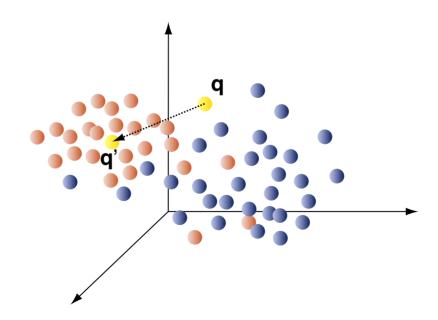


Query Refinement: Relevance Feedback

Given a result set R for a query q, and subsets $R^+ \subseteq R$ and $R^- \subseteq R$ of relevant and irrelevant documents, where $R^+ \cap R^- = \emptyset$, the query representation \mathbf{q} can be refined using Rocchio's update formula:

$$\mathbf{q}' = \alpha \cdot \mathbf{q} + \beta \cdot \frac{1}{|R^+|} \sum_{d^+ \in R^+} \mathbf{d}^+ - \gamma \cdot \frac{1}{|R^-|} \sum_{d^- \in R^-} \mathbf{d}^-,$$

where α , β , and γ adjust the impact of original query and (ir)relevant documents.



Query Refinement: Relevance Feedback

Given a result set R for a query q, and subsets $R^+ \subseteq R$ and $R^- \subseteq R$ of relevant and irrelevant documents, where $R^+ \cap R^- = \emptyset$, the query representation \mathbf{q} can be refined using Rocchio's update formula:

$$\mathbf{q}' = \alpha \cdot \mathbf{q} + \beta \cdot \frac{1}{|R^+|} \sum_{d^+ \in R^+} \mathbf{d}^+ - \gamma \cdot \frac{1}{|R^-|} \sum_{d^- \in R^-} \mathbf{d}^-,$$

where α , β , and γ adjust the impact of original query and (ir)relevant documents.

Observations:

- \Box Terms not in query q may get added; often a limit is imposed (say, 50).
- □ Terms may accrue negative weight; such weights are set to 0.
- Moves the guery vector closer to the centroid of relevant documents.
- Works well if relevant documents cluster; less suited for multi-faceted topics.

Relevance feedback can be obtained directly from the user, indirectly through user interaction, or automatically assuming the top-retrieved documents as relevant.

Discussion

Advantages:

- Severely improved retrieval performance compared to Boolean retrieval
- Partial query matching: not all query terms need be present in a document for it to be retrieved
- \lnot The relevance function ρ defines a ranking among retrieval documents with respect to their computed similarity to the query

Disadvantages:

Index terms are assumed to occur independent of one another