Vector space model

From Wikipedia, the free encyclopedia

Jump to: navigation, search

Vector space model (or *term vector model*) is an algebraic model for representing text documents (and any objects, in general) as <u>vectors</u> of identifiers, such as, for example, index terms. It is used in <u>information filtering</u>, <u>information retrieval</u>, <u>indexing</u> and relevancy rankings. Its first use was in the <u>SMART Information Retrieval System</u>.

Contents

[<u>hide</u>]

- <u>1 Definitions</u>
- <u>2 Applications</u>
- <u>3 Example: tf-idf weights</u>
- <u>4 Advantages</u>
- <u>5 Limitations</u>
- <u>6 Models based on and extending the vector space model</u>
 - 7 Software that implements the vector space model
 - <u>7.1 Free open source software</u>
- <u>8 Further reading</u>
- <u>9 See also</u>
- <u>10 References</u>

[edit] Definitions

Documents and queries are represented as vectors.

$$d_{j} = (w_{1,j}, w_{2,j}, \dots, w_{t,j})$$

$$q = (w_{1,q}, w_{2,q}, \dots, w_{t,q})$$

Each <u>dimension</u> corresponds to a separate term. If a term occurs in the document, its value in the vector is non-zero. Several different ways of computing these values, also known as (term) weights, have been developed. One of the best known schemes is <u>tf-idf</u> weighting (see the example below).

The definition of *term* depends on the application. Typically terms are single words, <u>keywords</u>, or longer phrases. If the words are chosen to be the terms, the dimensionality of the vector is the number of words in the vocabulary (the number of distinct words occurring in the <u>corpus</u>).

Vector operations can be used to compare documents with queries.

[edit] Applications



<u>Relevancy rankings</u> of documents in a keyword search can be calculated, using the assumptions of <u>document similarities</u> theory, by comparing the deviation of angles between each document vector and the original query vector where the query is represented as same kind of vector as the documents.

In practice, it is easier to calculate the <u>cosine</u> of the angle between the vectors instead of the angle:

A cosine value of zero means that the query and document vector are <u>orthogonal</u> and have no match (i.e. the query term does not exist in the document being considered). See <u>cosine similarity</u> for further information.

[edit] Example: tf-idf weights

In the classic vector space model proposed by <u>Salton</u>, Wong and Yang ^[1] the term specific weights in the document vectors are products of local and global parameters. The model is known as <u>term frequency-inverse document frequency</u> model. The weight vector for document *d* is , where

and

• tf_t is term frequency of term t in document d (a local parameter)

• is inverse document frequency (a global parameter). |D| is the total number of documents in the document set; is the number of documents containing the term *t*.

Using the cosine the similarity between document d_i and query q can be calculated as:

In a simpler <u>Term Count Model</u> the term specific weights do not include the global parameter. Instead the weights are just the counts of term occurrences: $w_{t,d} = tf_t$.

[edit] Advantages

The vector space model has the following advantages over the **Standard Boolean model**:

- 1. Simple model based on linear algebra
- 2. Term weights not binary
- 3. Allows computing a continuous degree of similarity between queries and documents
- 4. allows ranking documents according to their possible relevance
- 5. Allows partial matching

[edit] Limitations

The vector space model has the following limitations:

- 1. Long documents are poorly represented because they have poor similarity values (a small scalar product and a large dimensionality)
- 2. Search keywords must precisely match document terms; word <u>substrings</u> might result in a "<u>false positive</u> match"
- 3. Semantic sensitivity; documents with similar context but different term vocabulary won't be associated, resulting in a "<u>false negative</u> match".
- 4. The order in which the terms appear in the document is lost in the vector space representation.
- 5. Assumes terms are independent
- 6. Weighting is intuitive but not very formal

[edit] Models based on and extending the vector space model

Models based on and extending the vector space model include:

- <u>Generalized vector space model</u>
- (enhanced) Topic-based Vector Space Model
- Latent semantic analysis
- Latent semantic indexing
- DSIR model
- <u>Term Discrimination</u>
- <u>Rocchio Classification</u>

[edit] Software that implements the vector space model

The following software packages may be of interest to those wishing to experiment with vector models and implement search services based upon them.

[edit] Free open source software

- <u>Apache Lucene</u>. Apache Lucene is a high-performance, full-featured text search engine library written entirely in Java.
- <u>SemanticVectors</u>. Semantic Vector indexes, created by applying a Random Projection algorithm (similar to <u>Latent semantic analysis</u>) to term-document matrices created using Apache Lucene.

[edit] Further reading

- <u>G. Salton</u>, A. Wong, and C. S. Yang (1975), "<u>A Vector Space Model for Automatic Indexing</u>," *Communications of the ACM*, vol. 18, nr. 11, pages 613–620. *(The article in which the vector space model was first presented)*
- Description of the vector space model
- Description of the classic vector space model by Dr E Garcia

[edit] See also

- Inverted index
- <u>Compound term processing</u>

[edit] References

1. <u>^ G. Salton , A. Wong , C. S. Yang, A vector space model for automatic indexing</u>, Communications of the ACM, v.18 n.11, p.613-620, Nov. 1975

Retrieved from "<u>http://en.wikipedia.org/wiki/Vector_space_model</u>" Categories: Information retrieval