# Vector Space Model

The vector space model is a mathematical model used in information retrieval (IR) to represent text documents and queries as vectors in a high-dimensional space. It is based on the idea that the meaning of a document or a query can be inferred from the terms that occur in it.

In the vector space model, each term in the document or query is represented by a dimension in the vector space, and the value in that dimension corresponds to the frequency of the term in the document or query. Thus, each document or query can be represented as a vector of term frequencies, also known as a bag-of-words model.

The similarity between two documents or a document and a query can be measured using various similarity measures, such as cosine similarity, Jaccard similarity, and Euclidean distance. The most common measure used is cosine similarity, which computes the cosine of the angle between two vectors. The closer the cosine similarity between two vectors is to 1, the more similar the documents or query are.

The vector space model is widely used in IR systems such as search engines, where it is used to retrieve relevant documents based on user queries. It has also been applied in other fields, such as natural language processing, text classification, and sentiment analysis.

## Objective Vector Space Model:

The objective of the vector space model is to represent text documents and queries as vectors in a high-dimensional space, where each dimension corresponds to a term in the document or query. The main goal of this model is to capture the meaning and context of the text by representing it in a mathematical form that can be processed by algorithms.

The vector space model is used primarily in information retrieval to rank documents based on their relevance to a query. The model assumes that the relevance of a document to a query can be estimated based on the similarity between their vector representations in the high-dimensional space. Thus, the objective of the vector space model is to enable efficient and accurate retrieval of relevant documents from a large collection based on a user's query.

Another objective of the vector space model is to enable text analysis and informationextraction from large volumes of text data. By representing text in a mathematical form, it becomes

possible to apply various algorithms and techniques from machine learning and natural language processing to analyze the data and extract insights from it.

**How Vector Space Model Work?**

The vector space model provides a mathematical framework for representing text documents and queries as vectors in a high-dimensional space. The functionality of the vector space model can be summarized as follows:

1. Document and query representation: The vector space model represents each document and query as a vector of term frequencies, where each dimension corresponds to a term in the document or query. This representation captures the important terms and their frequencies in the text, enabling efficient and accurate retrieval of relevant documents from a large collection.

2. Term weighting: The vector space model uses term weighting to give more weight to importantterms and less weight to common terms that may not be as relevant to the document or query.There are several term weighting schemes, such as TF-IDF, which measures the importance ofa term in a document based on its frequency and rarity in the collection.

3. Similarity measure: The vector space model uses a similarity measure, such as cosine similarityor Jaccard similarity, to measure the similarity between two vectors in the high-dimensional space. This enables the model to retrieve documents that are most similar to a query based on their vector representations.

4. Query expansion: The vector space model can be used to expand a query by adding synonymsor related terms to it. This is done by finding terms that are semantically similar to the originalquery terms based on their vector representations in the high-dimensional space.

5. Dimensionality reduction: The vector space model can be used to reduce the dimensionality ofthe high-dimensional space by applying techniques such as singular value decomposition (SVD) or principal component analysis (PCA). This can improve the efficiency and accuracy of the model by reducing the number of dimensions while preserving the most important information.

Overall, the functionality of the vector space model enables efficient and accurate retrieval of relevant documents from a large collection based on a user's query, as well as text analysis and

information extraction from large volumes of text data.

**Advantages of Vector Space Model:**

1. Flexibility: The vector space model is flexible and can be used for a wide range of text analysistasks, such as information retrieval, text classification, and sentiment analysis.

2. Efficient retrieval: The vector space model allows for efficient retrieval of relevant documentsfrom a large collection based on a user's query, using similarity measures such as cosine similarity.

3. Term weighting: The vector space model uses term weighting to give more weight to importantterms and less weight to common terms, improving the accuracy of retrieval and analysis.

4. Query expansion: The vector space model can be used to expand a query by adding synonymsor related terms, improving the precision and recall of retrieval.

5. Text analysis: The vector space model enables text analysis and information extraction from large volumes of text data, using techniques from machine learning and natural language processing.

**Disadvantages of Vector Space Model:**

1. Sparse representation: The vector space model results in a sparse representation of text data, with many dimensions having zero values, which can lead to computational complexity and memory issues.

2. Over-reliance on term frequency: The vector space model relies heavily on term frequency to represent text data, which can result in inaccuracies when terms have different meanings in different contexts.

3. Synonymy and polysemy: The vector space model can have difficulties dealing with synonymy (different terms with the same meaning) and polysemy (the same term with multiple meanings),which can result in inaccurate retrieval and analysis.

4. Dimensionality: The vector space model can have high dimensionality, especially for large

collections of text data, which can lead to computational complexity and overfitting.
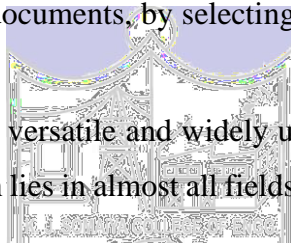
5. Lack of context: The vector space model does not capture the context of the text, which can result in inaccuracies when dealing with ambiguous or complex text data

**Applications of Vector Space Model**

The vector space model has various applications in natural language processing, information retrieval, and text analysis, including:

1. Information Retrieval: The vector space model is widely used in information retrieval systems,such as search engines, to retrieve relevant documents from a large collection based on a user'squery.

2. Text Classification: The vector space model is used in text classification tasks, such as sentiment analysis and topic modeling, to classify text into different categories based on their vector representations.

3. Question Answering: The vector space model is used in question answering systems to retrieveanswers from a large knowledge base based on the similarity between the question and the documents.

4. Recommender Systems: The vector space model is used in recommender systems to recommend products or services to users based on their preferences, by finding similar items based on their vector representations.

5. Clustering: The vector space model is used in clustering algorithms to group similar documents together based on their vector representations, which can be used for text segmentation or topicmodeling.

6. Machine Translation: The vector space model is used in machine translation systems to align and translate text from one language to another, by representing text in a high-dimensional space and finding the most similar translation.

7. Text Summarization: The vector space model is used in text summarization systems to generatesummaries of long documents, by selecting the most important sentences based on their vectorrepresentations.

Overall, the vector space model is a versatile and widely used tool in natural language processing and text analysis, with application lies in almost all fields including information retrieval, text

classification, question answering, recommender systems, clustering, machine translation, andtext summarization.

**Example of Vector Space Model**

document 1 = "The cat sat onthe mat cat mat"

document 2 = "The cat perched on the mat."

document 3 = "The feline rested on the mat."

Query = "cat AND mat NOT sat"

The term-document incidence matrix

|  | Document 1 | Document 2 | Document 3 |
|---|---|---|---|
| Cat | 1 | 1 | 0 |
| Feline | 0 | 0 | 1 |
| Mat | 1 | 1 | 1 |
| On | 1 | 1 | 1 |
| Perched | 0 | 1 | 0 |
| Rested | 0 | 0 | 1 |
| Sat | 1 | 0 | 0 |
|  |  |  |  |

| The | 1 | 1 | 1 |
|---|---|---|---|

term document frequency table with "NOT sat"

|  | Document 1 | Document 2 | Document 3 |
|---|---|---|---|
| Cat | 1 | 1 | 0 |
| Feline | 0 | 0 | 1 |
| Mat | 1 | 1 | 1 |
| On | 1 | 1 | 1 |
| Perched | 0 | 1 | 0 |
| Rested | 0 | 0 | 1 |

| NOT sat | 0 | 1 | 1 |
|---------|---|---|---|
| The | 1 | 1 | 1 |

document 1 : 1 A 0A 0 = 0 document

2 : 1 A 1 A 1 = 1 document 3 : 10A 1

A 1 = 0

Hence, the most relevant is document 2

# Boolean Model

The first search engines used the Boolean retrieval model, and it is still used by search engines today.It is also called "exact-match retrieval" because documents are only retrieved if they match the query exactly. If they don't, they are not retrieved.The name "Boolean" comes from the fact that there are only two possible answers to a query (TRUE or FALSE) and that most queries are written using operators from Boolean logic (AND, OR, NOT).

**Boolean Model**

The objective of the boolean retrieval model is to retrieve relevant documents from a collection of documents based on the presence or absence of specific keywords or combinations of keywords in the document. The model is based on the principle of boolean logic, which uses logical operators such as AND, OR, and NOT to combine keywords and refine search results.

The boolean retrieval model aims to provide precise search results by allowing users to specify exactly which keywords or combination of keywords they are looking for. By using boolean operators, users can refine their search queries and narrow down the results to the most relevant documents.

Another objective of the boolean retrieval model is to provide a simple and efficient method for information retrieval. The model is easy to implement and can handle large volumes of data efficiently, making it a popular choice for search engines and other information retrieval systems.

Overall, the objective of the boolean retrieval model is to provide users with a powerful tool for retrieving relevant documents quickly and accurately. By using boolean logic and indexing techniques, the model can efficiently search through large collections of documents and retrieve only those that match the user's query.

## **Working of Boolean Model Works**

The boolean retrieval model is a classical information retrieval model that uses Boolean logic to retrieve relevant documents from a collection of documents based on the presence or absence of specific keywords or combinations of keywords in the document. The model works on the principle of boolean logic, where the documents are represented as sets of words or terms, and queries are represented as boolean expressions.

The model performs the following functions:

1. Indexing: The boolean retrieval model first indexes the collection of documents by creating aninverted index, which is a data structure that maps each term or keyword in the collection to alist of documents that contain that term.

2. Query processing: When a user enters a query, the boolean retrieval model processes the queryby translating it into a boolean expression, which is a combination of terms and Boolean operators (AND, OR, NOT). The model then applies the boolean operators to the inverted indexto retrieve the relevant documents that match the query.

3. Ranking: The boolean retrieval model <u>does not rank </u>the retrieved documents based on their relevance to the query. Instead, it returns all the documents that match the query. It is up to theuser to determine which documents are most relevant.

4. Precision and recall: The boolean retrieval model calculates the precision and recall of the retrieved documents. Precision is the percentage of relevant documents retrieved out of the total number of documents retrieved, while recall is the percentage of relevant documents retrieved out of the total number of relevant documents in the collection.

Overall, the boolean retrieval model is a simple and efficient method for retrieving relevant documents based on keyword queries. Its functionality is based on boolean logic and it is widely used in information retrieval systems such as search engines and databases.

**Advantages of Boolean Model**

1. Clean and Easy to implement - Intuitive concept
2. If the resulting document set is either too small or too big, it is directly clear which operatorswill produce respectively a bigger or smaller set.

3. It gives (expert) users a sense of control over the system. It is immediately clear why adocument has been retrieved given a query.

### Disadvantages of Boolean Model

1. Exact matching may retrieve too few or too many documentsHard to translate a query into a Boolean expression

2. All terms are equally weighted

3. More like data retrieval than information retrieval

4. Retrieval based on binary decision criteria with no notion of partial matchingNo ranking of the documents is provided (absence of a grading scale)

5. Information need has to be translated into a Boolean expression, which most users find awkward.

6. The Boolean queries formulated by the users are most often too simplistic

7. The model frequently returns either too few or too many documents in response to a user query

### Application of Boolean Model

1. Information Retrieval: The boolean retrieval model is widely used in search engines to retrieve relevant documents based on specific keywords or phrases. It is an effective way tofilter irrelevant documents and present only the most relevant ones to the user.

2. Database Management: The boolean retrieval model can be used to retrieve data from a database based on specific criteria. For example, if a company wants to retrieve all the salesdata for a specific product, they can use boolean queries to retrieve the relevant data.

3. Spam Filtering: The boolean retrieval model is used in spam filters to identify and filter outunwanted emails based on specific keywords or phrases that are commonly found in spam emails.

4. Legal Research: The boolean retrieval model is used in legal research to retrieve relevant caselaw or statutes based on specific keywords or phrases.

5. Medical Research: The boolean retrieval model is used in medical research to retrieve relevantarticles or studies based on specific medical terms or conditions.

6. Overall, the boolean retrieval model is a powerful tool for information retrieval and is used ina wide range of applications where precise searching and filtering of data is necessary.

**Illustration with Example:**

We will be using a tf-idf weighting scheme to represent the documents as vectors. Other popular methods under vector space models are bag of words models, word2vec etc.

Example:

Let us take following documents

d1: 'The cat sat on the mat cat mat'
d2: 'The cat perched on the mat' d3:
'The feline rested on the mat'

Let us calculate tf-idf for word 'cat' for all the documents.

$tf('cat', d1) = 0.04$

$tf('cat', d2) = 0.0068$

Using $idf tf(('cat 'cat above', ') d$ 3values=) $log=$ tf-idf0(.32 0) can be 0.176

calculated

as:

$tf - idf('cat', d1, D) = 0.0704$

$tf - idf('cat', d2, D) = 0.0587$

$tf - idf('cat', d3, D) =$ (Autonomous 0.0

Using the above steps, we generate the scores for all the words in the corpus.

KJSCE/IT/TYBTech/SEMVIII/IR/2022-23

|      | cat    | sat    | mat | rested | feline | perched |
|------|--------|--------|-----|--------|--------|---------|
| doc1 | 0.0704 | 0.0954 | 0.0 | 0.000  | 0.000  | 0.000   |
| doc2 | 0.0587 | 0.0000 | 0.0 | 0.000  | 0.000  | 0.159   |
| doc3 | 0.0000 | 0.0000 | 0.0 | 0.159  | 0.159  | 0.000   |

**Vector representation of documents:**

vocabulary :      {'cat', 'feline', 'mat', 'perched', 'rested',

'sat'} doc1:      {0.0, 0.0, 0.0704, 0.0, 0.0, 0.0954} doc2:

{0.0, 0.0, 0.0587, 0.159, 0.0, 0.0} doc3:   {0.0, 0.159, 0.0,

0.0, 0.159, 0.0} For query q: cat mat

vector representation: {1, 0, 1, 0, 0, 0}

Ranking is obtained using cosine

similarity:

sim(q, d1) = 0.3659 sim(q, d2) = 0.2449 sim(q, d3) = 0

Thus retrieval ranking would be d1, d2 and d3. d3 being completely irrelevant to the queryas similarity score is 0. d1>d2>d3

**Key Differences Between Boolean and Vector Model:**

| IR models(IR mod)/ attributes(A) | Boolean (IR mod) | Vector space (IR mod) |
|---|---|---|
| Concept(A) | Based on set theory and Boolean algebra | Based on the concept of vectors |
| Representation(A) | Documents are represented by the index terms extracted from documents, and queries are Boolean expressions on terms. | Represented in the form of weighted-term vectors. Cosine measure is used to find the similarities |
| Information type(A) | Does not consider semantic information | It consider semantic information |
| Word occurrence(A) | Number of occurrence are not mentioned | Tells about the number of occurrence |
| Output(A) | Exact match of the output to the query | Best match of the query |
| Advantages(A) | Easy to implement | Simple model, weights are not in binary |
| Disadvantages(A) | Does not rank documents, retrieves too many or too few | Suffers from synonymy and polysemy. It theoretically assumes that terms are statistically independent |