

Dated:

CHAPTER #6: VECTOR SPACE Model

- ★ Metadata => Digital docs generally encode in machine-recognizable form, certain metadata like authors, title and date of publication.
- ★ field => The metadata values are called fields.
- ★ Parametric index => The system holds the inverted index as well as a parametric index which exists for all fields. These search parameters retrieve information and then are merged with inverted index for final results.
e.g: date of creation parameter can be used to get docs after a certain date. The dict comes from a fixed vocab.
- ★ zone => A zone is an arbitrary unbounded amount of text like abstracts or references in research papers. We may build a separate inverted index for each zone to support queries. The dict must structure whatever vocab stems from the zone.

Author: _____
title: _____
language: _____ } Parametric search example

William.abstract → [1] → [121] → [1441] } zone index

William.title → [2] → [4] → [8] } zones as ext of dict

William → [2.author, 2.title] → [3.author] } zone is encoded
in postings.

★ Weighted Page Ranking

★ Ranked Retrieval =>

↳ return an ordering over the top documents in the collection.

↳ the user's query is in human language.

X

- * Score-ranking =>
 - ↳ for ranked retrieval, assign a score $[0, 1]$ - to each document.
 - ↳ this score measures how well doc & query match.
 - ↳ if score = 0, query is not in document.
 - ↳ the higher the score, the more frequent the term occurs in doc.
 - ↳ the most basic is the Jaccard co-efficient.

- * Weighted Zone Scoring => Different zones may have a different importance in evaluating how a document matches a query.

Given a boolean query q_j and a document d , weighted zone scoring assigns to the pair (q_j, d) a score in the interval $[0, 1]$ by computing a linear combination of zone scores.

Consider a set of documents each of which has l zones. Let $q_1, \dots, q_l \in [0, 1]$ such that $\sum_{i=1}^l q_i = 1$.
 zones of a doc. \rightarrow only 0 or 1

For $1 \leq i \leq l$, let s_i be the boolean score denoting a match (or mismatch) b/w q_j and i th zone, then the weighted zone scoring becomes:

$$\sum_{i=1}^l q_i s_i$$

It is also called ranked boolean retrieval.

↳ Example: Consider the query Shakespeare in a collection which in each document has zones: author, title, body. The Boolean score will be 1 if query exists in zone else 0. Let $q_1 = 0.2$ by author weight, $q_2 = 0.3$ be title weight, $q_3 = 0.5$ be body weight. This suggest that author zone is least important and body zone is most important.

Thus if, term was not to appear in author, the document score will be 0.8.

* term frequency (tf) \Rightarrow number of occurrences of a term in the document.

* Bag of words \Rightarrow In this model, a text is represented as the bag(multiset) of its words, disregarding grammar and even word order but keeping multiplicity.

ex: John likes to watch movies. Mary likes movies too.

BoW = {John, like, watch, movie, Mary, like, movie}.

To compute scores, we find BoW of the query & find the jaccard co-efficient.

↳ Issues =>

↳ Doc. size affects overall similarity.

↳ frequency of common terms are ignored.

↳ similarity is based on common terms, which are treated independent of each other.

* Document Frequency (df) => number of docs in the collection that contain the term t.

* Inverse Doc. Freq. (idf) =>

Let N be the number of documents in a collection, the idf for a term t is:

$$\text{idf}_t = \log \frac{N}{\text{df}_t}$$

The idf of a rare term is high, and of a frequent term is low.

* tf-idf weighting => we combine the two to produce a composite weight for each term in a document.

The tf-idf weighting scheme assigns to a term t a weight in the document d as:

$$\text{tf-idf}_{t,d} = \text{tf}_{t,d} \times \text{idf}_t$$

↳ highest when t occurs many times within a small number of documents. (lending high discriminating power)

↳ lower when term occurs fewer times in a document, or occurs in many documents. (less pronounced relevance signal)

↳ lowest when the term occurs in virtually all documents.

* Cumulative frequency => show many times a term appears in a collection (cf).

* Document Vectors =>

At this point, viewing each document has a vector with one component corresponding to each term in the dictionary, together with the weight. If they do not occur in the doc, the weight is zero.

↳ Overlap score measure => the score of a document d is the sum over all query terms, of the number of times each query terms occur in d . We can refine this with the use of tf-idf.

$$\text{Score}(q, d) = \sum_{t \in q} \text{tf-idf}_{t,d}$$

* Vector Space Model => The representation of a set of documents as vectors in a common vector space and it is fundamental to a host of info retrieval operations ranging from scoring docs on a query, document classification & document clustering.

Basic Properties

↳ Idea =>

↳ The basic idea of VSM is to have a high-dimensional hyperspace where each term is a dimension and document & query can be represented in this space.

↳ It is a simple yet effective method of designing ranking functions for IR.

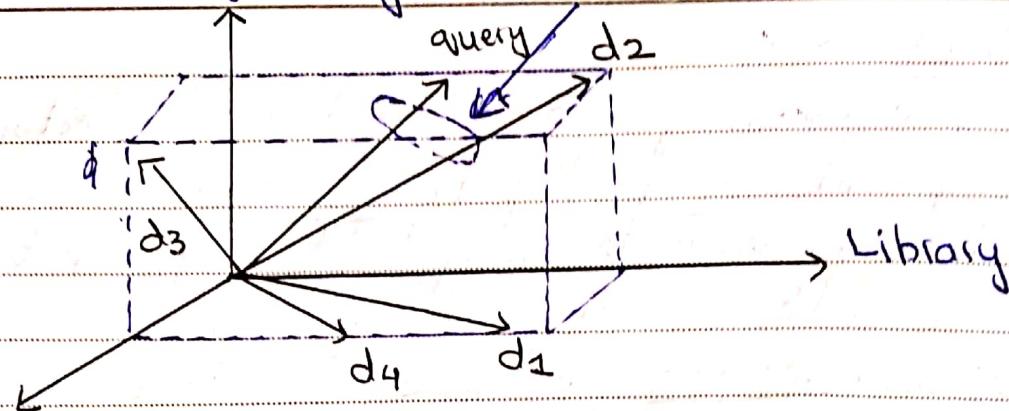
↳ Assumptions =>

↳ All terms are features and they are all independent.

↳ the vocabulary (V) is limited.

↳ the vocabulary formed a space of $|V|$ dimensions.

↳ Relevance can be computed as distance or angle b/w vectors of query and document.



Presidential

→ Distances Measurements =>

Let $\vec{v}(d)$ denote the vector derived from document d . How do we quantify the similarity b/w two documents in the vector space.

The first step attempt might be to consider the magnitude of vector difference b/w two vectors. This suffers from a major drawback if the documents have very similar content but one of the documents is too long. Thus the relative distribution of terms may be identical but absolute term frequencies of one may be far larger.

To compensate for this effect, the standard way of quantifying the similarity is to compute the Cosine similarity.

$$\text{sim}(d_1, d_2) = \frac{\vec{v}(d_1) \cdot \vec{v}(d_2)}{\|\vec{v}(d_1)\| \|\vec{v}(d_2)\|}$$

where the numerator is the dot product while the denominator is the Euclidean length. The dot product of two vectors x and y can be calculated as $\sum_{i=1}^M x_i y_i$, where M is the number of components and the euclidean length is $\sqrt{\sum_{i=1}^M \vec{v}_i^2(d)}$.

The effect of denominators is to length normalize the vectors to unit vectors $\vec{u}(d_1) = \vec{v}(d_1) / \|\vec{v}(d_1)\|$ and $\vec{u}(d_2) = \vec{v}(d_2) / \|\vec{v}(d_2)\|$.

$$\cos \theta = \frac{\mathbf{q}_1 \cdot \mathbf{d}}{\|\mathbf{q}_1\| \|\mathbf{d}\|}$$

* The component magnitudes are tf instead of exist or not exists
Dated: _____

So, the eqn becomes

$$\text{sim} \leftarrow \text{sim}(d_1, d_2) = \vec{v}(d_1) \cdot \vec{v}(d_2)$$

The term documents can be represented with a term document matrix of MxN dimensions where M is the total terms and N is the documents in the collection.

↳ Queries as vectors =>

We are representing queries as vectors as well. The key idea is to assign each document d a score equal to dot product $\Rightarrow \vec{v}(q_1) \cdot \vec{v}(d)$ or

by viewing query as a bag of words, we are able to treat it as a very short document. Hence, we can use cosine similarity to find that document score:

$$\text{Score}(q_1, d) = \frac{\vec{v}(q_1) \cdot \vec{v}(d)}{\|\vec{v}(q_1)\| \|\vec{v}(d)\|}$$

Computing the cosine similarities bw query and each document, sorting the resulting scores and selecting the top K documents is really expensive.

↳ Instantiation of XSM =>

↳ XSM is a framework

↳ Computing Vector Scores =>

COSINE SCORE(q_1):

- 1 float Scores[N] $\leftarrow 0$
- 2 Initialize length[N]
- 3 for each query term t
- 4 calculate W_{t,q_1} and fetch posting list for t
- 5 for each pair $(d, f_{t,d})$ in postings list
- 6 add $wf_{t,d}$ to Scores[d]
- 7 Read the array length[d]
- 8 for each d
- 9 divides Scores[d] by length[d]
- 10 return top K scores.

We are seeking the K documents with the highest vector space scores on the given query.

The array lengths[N] holds the lengths(normalization factors) of each of the N documents. The array Scores holds the document scores. The outermost loop in Step 3 repeats the updating of scores iterating over each query term. In Step 5, we calculate the weight in query vector for term t . Step 6-8 updates the score of each document by adding the contribution from t . This process of adding in contributions one query term at a time is known as term-at-a-time scoring or accumulation, and the N elements of Scores array are known as accumulators.

For this purpose, it seems necessary to store the weight of term with each posting entry. This is wasteful because it requires storing an extra float.

Two ideas can resolve this problem.

First, if we are using ~~inter~~ idf, we need not precompute idf, it suffices to store N/dft at the head of the postings for t . Second, we store the term frequency for each posting entry. Finally, a priority queue with a max heap taking $2N$ comparisons to construct and $O(\log n)$ extraction can be used to get the top K scores.

\hookrightarrow document-at-a-time \Rightarrow computing scores of one doc at a time.

\hookrightarrow advantages \Rightarrow

\hookrightarrow simple model based on linear algebra.

\hookrightarrow Term weights are not binary.

\hookrightarrow allows computing a continuous degree of similarity.

\hookrightarrow allows ranking documents.

\hookrightarrow allows partial matching.

↳ disadvantages =>

↳ the order in which terms appear is lost.

↳ theoretically assumes that terms are independent.

↳ search keywords must precisely match doc terms,
word substrings might give false res.

↳ semantic sensitivity; documents with similar context

but different term vocab won't be associated resulting
in a false -ve.

* Sublinear tf scaling => {reduce impact of high frequency counts}

A common modification can be made to normal tf since

all occurrences may or may not signify the same importance

so the new weight can be calculated as:

$$W_{tf,d} = \begin{cases} 1 + \log tf_{t,d} & \text{if } tf > 0 \\ 0 & \text{otherwise} \end{cases}$$

tf-idf can also be changed as:

$$W_{tf-idf,t,d} = W_{tf,t,d} \times idf_t$$

* Maximum tf normalizations

tf-idf scores sometimes give skewed results based on frequency so this issue can be addressed by dividing the raw frequency count of each term by the highest frequency count of any term in the document. This ensures that the highest score can be 1 and all others are proportionally adjusted.

$$n_{tf,t,d} = a + (1-a) \frac{tf_{t,d}}{tf_{max}(d)}$$

Here, a is called a smoothing term that has a value b/w 0 and 1 usually 0.4 is used for damping the contribution of the second term which is scaling down tf by the largest tf value in d.

↳ unstable changes in stopwords can cause drastic changes.

↳ a doc may have an unusually frequent outlier.