

# Information retrieval

Francesco Tomaselli

March 19, 2021

## Contents

<b>1</b>	<b>Vector space model</b>	<b>2</b>
1.1	Tokenization . . . . .	2
1.2	Normalization . . . . .	3
1.3	Term Weighting . . . . .	3
<b>2</b>	<b>Evaluation on retrieval systems</b>	<b>4</b>
2.1	The notion of quality . . . . .	4
2.2	Quality measures . . . . .	5
2.3	Evaluation of ranking systems . . . . .	6
<b>3</b>	<b>Relevance feedback and query expansion</b>	<b>7</b>
3.1	Local methods . . . . .	7
3.2	Global Methods . . . . .	8
3.2.1	How to add information? . . . . .	8
3.2.2	Spelling correction . . . . .	9

# 1 Vector space model

**Text transformation** The first step to process queries on text is to transform it into a model that can be computed by a machine.

Such a model should be suitable for any kind of text, should be able to capture the semantic of words and then should support the notion of *text similarity*.

The main idea of a vector space model is to map documents in a vector, so they appear as points in a certain space.

For instance, we can collect some words and index them, then build a vector for each document

$$d_k = [x_0, \dots, x_n]$$

where the value  $x_i$  stores the occurrences of the word indexed with  $i$  in that document.

**Similarity** This model somehow implements the notion of text similarity, represented by some sort of distance function between two documents in space, perhaps the *Euclidean distance* or the difference between the angle of the two vectors. There are a lot of distance measures, and they vary in applications. Some of them assume normalization in the document vectors.

**Document-Term matrix** If we arrange document vectors in a vector we obtain the *document-term matrix*, or the *term-document matrix* if we take the opposite. The number of dimensions in the space is equal to the number of words in the vocabulary.

That can be a problem without some sort of lemmatisation and stemming, as sparsity and dimension of the matrix increases.

**Vector length problem** By considering the length of the vector in a similarity measure, we are introducing some problems, let's just consider a text and its abstract, the Euclidean distance would be huge, as most of the words do not appear that often.

To solve the problem we can consider the *Cosine similarity* or normalize vector lengths.

The main phases in the generation of a vector space models are: tokenization, normalization, weight, and indexing.

## 1.1 Tokenization

The goal is to split the text in words, for instance we could split the document at the spaces.

In practice this can be obtained with:

- *Regex*
- *Corpus* based
- *Machine Learning* based

## 1.2 Normalization

The idea is to normalize each token, that could mean writing verbs in the normal form, or plurals at the singular.

**Stemming** One form of normalization is stemming, that simply truncate words to obtain singular or infinite forms. The problems here could be creation of meaningless words and also close words in meaning could be completely different. An example is the *Porter* stemming algorithm.

**Lemmatisation** Another way to normalize is to lemmatise, It's similar to stemming, but the words are replaced with the dictionary root, this also has a problem, indeed sometime mapping of same words to different lemmas can happen. A way to lemmatise is to use a *Dictionary based* lemmatisation, or something more complicated such as *Wordnet*.

**Wordnet** Wordnet is a large lexical database, verbs, nouns, adjectives and adverbs are grouped into cognitive synonyms, called synsets, each expressing a different concept. Synsets are connected to each other by conceptual relations and they are also connected to a lemma.

## 1.3 Term Weighting

Weighting the words is a crucial point in text processing, an easy way could be counting the frequency of the terms in the document.

**Term frequency** An example of weighing is the term frequency, which counts the frequency of a term given a document:

- *Boolean*: so a 0 or 1 if the word is present or not
- *Natural*:  $tf_{t,d}$  = count of a specific word
- *Log*:  $1 + \log(tf_{t,d})$
- *Augmented*:  $0.5 + \frac{0.5tf_{t,d}}{\max_{t'} tf_{t',d}}$
- *Log average*:  $0.5 + \frac{0.5 \log(tf_{t,d})}{1 + \log(\max_{t'} tf_{t',d})}$
- *Max tf norm*:  $k + (1 - k) \frac{tf_{t,d}}{tf_{\max}(d)}$

**Stop words problem** A problem in this model is the excessive presence of stop words and punctuations. To compensate, we can remove them in the normalization phase.

But sometimes the stop words can be important, for instance with phrasal verbs. A solution can be to count the number of documents that contains that given word.

**Inverse document frequency** The idea is to count the number of documents that contains a given words, and to prefer less frequent words:

$$idf = \log \frac{N}{df_t}$$

There are many variants:

- *Smooth*:  $\log(\frac{N}{1+df_t}) + 1$
- *Max*:  $\log(\frac{\max_{t' \in d} df_{t'}}{1+df_t}) + 1$
- *Probabilistic*:  $\log \frac{N-df_t}{df_t}$

**TfIdf** By multiplying term frequency and inverse term frequency we obtain this metric.

$$TfIdf(t, d) = tf_{t,d} \cdot idf_t$$

Regarding the value of terms:

- the terms with a high *TfIdf* usually appears in a small number of documents
- the metric is low when a term appears a few times in a document or when it appears in many documents

## 2 Evaluation on retrieval systems

The goal of evaluation is to assess the quality of the results obtained by an IR system. There's the need of knowing a *ground truth*, so an annotated corpus where for each task we know what documents are relevant. The annotations could be created manually or derived from the data, if it contains annotations.

### 2.1 The notion of quality

Given a corpus  $C$  and a query  $q$ , the task is to find a set of documents  $A_{q,C}$  that match  $q$ , but after retrieving such documents, there's the need to estimate the quality of results.

**Precision** To formalise the quality of the retrieved documents  $A_{q,C}$  we could count how many of these documents are relevant to  $q$ , this is the notion of precision:

$$Prec = \frac{\text{relevant retrieved}}{\text{retrieved}}$$

Note that in order to know if a document is relevant or not we need the ground truth or a user feedback.

This measure does not suffice, as for instance, if we retrieve only one correct document we would have maximum precision.

**Recall** The precision measure does not take in consideration how many relevant documents are there, that is crucial to assess quality of a query result. So we can consider another measure, called Recall:

$$Rec = \frac{\text{relevant retrieved}}{\text{relevant}}$$

**Information need** Given the two quality measures, should we aim at better precision or more recall? Trivially both, but it actually depends on the information need.

In some cases a really high recall is not necessary, while other times a lower precision could be accepted while not missing anything relevant.

**F1 Score** To take into consideration both precision and recall, we could take a weighted mean, so a tradeoff between the two

$$F1 = \frac{2 \cdot Prec \cdot Rec}{Prec + Rec}$$

**Baseline system** To assess the quality of an IR system, we need a baseline system, something trivial such as tossing a coin to decide whereas a document is relevant or not. Given its quality measure, we can infer the quality of our, hopefully, more complex system, in fact, quality measures can't be seen as absolute values, there's always the need to compare.

## 2.2 Quality measures

To formally define the notions introduced in the previous section, we need to take into consideration a more detail measure for errors.

**Confusion Matrix** Given a query  $q$ , a ground truth  $E_q$  of relevant documents with respect to  $q$ , and a set of retrieved documents  $A_q$ , we define this matrix:

	$d \in E_q$	$d \notin E_q$
$d \in A_q$	True positive	False positive
$d \notin A_q$	False Negative	True Negative

So now we can redefine the *precision*, *recall* and *F1* measures as follows

$$Prec = \frac{TP}{TP + FP} \quad Rec = \frac{TP}{TP + FN} \quad F1 = \frac{TP}{TP + \frac{1}{2}(FP + FN)}$$

The actual confusion matrix is obtained by counting the documents given the retrieved ones and the ground truth.

The matrix can be useful to estimate the system parameters, for instance, if the number of retrieved documents is set to a certain value  $k$  and the value of true positives is really high, probably a smaller  $k$  would do the job, while if a lot positives are left out, i.e, the matrix has a high number of false negatives, maybe an higher  $k$  would make the system better.

**Other measures** The are a lot of measures to estimate the quality of a system, such as *specificity*, *negative predictive value*, *miss-rate*, *fall-out* and *accuracy*. Finding the right measure is really important, for instance, using accuracy with an unbalanced dataset is not better than precision, recall and F1.

## 2.3 Evaluation of ranking systems

We talked about boolean retrieval systems and we took into consideration if a document is retrieved or not.

In a raking scenario, we want to give importance not only to precision and recall, but also to the rank assigned by the system.

**Setting a threshold** One way to go is to set a specific threshold and compute the precision and recall of those top documents. This method does not take into consideration the ranking of the documents.

**Precision at K** If we take the first  $K$  retrieved documents, ordered by rank, we can estimate the quality of the rank by computing the precision for the top  $K$  documents. By iterating the threshold we can better estimate the performances, opposed to setting a unique threshold.

To decide what thresholds to use we could use the distribution of the system ranking.

**Discounted cumulative gain** This approach takes into consideration the ranking and discount the gain given by a document with respect to its position in it.

$$DCG = \sum_{i=1}^n \frac{R_j}{\log(i+1)}$$

where  $R_j$  is the rank assigned to a document.

**Precision vs Recall curve** Another approach is to compute precision and recall measures at each point in the ranking. If we take the measurements and plot them, we find a correlation between precision and recall and if we compute the integral of that function, we obtain the average precision.

To have a smoother curve, it's possible to interpolate the precision, so for each point in the recall axis, we take the maximum precision of consecutive points.

### 3 Relevance feedback and query expansion

The idea behind this section is the possibility to use some sort of feedbacks to tweak queries before passing them to the information retrieval system.

**Relevance feedback** It's any feedback we can collect about correctness of a certain query. For instance a user feedback or a indirect feedback about usefulness of an answer, i.e. the user clicked on a particular website

**Query expansion** Transform user queries exploiting the relevance feedback to improve system results over that query.

When expanding a query there are two possible methods, local and global methods, depending on whereas the expansion is applied to a single query or to all queries to the system:

- *Local methods*: related to a particular query, they are based on relevance and indirect feedback of that particular query
- *Global methods*: methods applied to all queries, those can be expansion with a thesaurus, or spelling correction

#### 3.1 Local methods

**Vector shifting** Given the sets  $R$  and  $N$  of relevant and non relevant results respectively, the task is to find a query vector  $\vec{q}_o$  that maximizes similarity with  $R$  and minimizes the one with  $N$ :

$$\vec{q}_o = \operatorname{argmax}_{\vec{q}} [\operatorname{sim}(\vec{q}, R) - \operatorname{sim}(\vec{q}, N)]$$

Given a similarity function, for instance the cosine similarity, the equation becomes

$$\vec{q}_o = \frac{1}{|R|} \sum_{\vec{d}_i \in R} \vec{d}_i - \frac{1}{|N|} \sum_{\vec{d}_j \in N} \vec{d}_j$$

where the two elements are the centroids of  $R$  and  $N$ .

**Rocchio algorithm** This approach is a generalization of the main idea introduced in the section, it proposes a new vector derived by making the query vector closer to the centroid of relevant documents.

$$\vec{q}_m = \alpha \vec{q} + \beta \frac{1}{|R|} \sum_{\vec{d}_i \in R} \vec{d}_i - \gamma \frac{1}{|N|} \sum_{\vec{d}_j \in N} \vec{d}_j$$

where  $\alpha$ ,  $\beta$  and  $\gamma$  can balance the role of each component.

One of the problems of this approach is the difficulty in estimating a document relevance, also, the sets  $R$  and  $N$  could be somehow merged.

**Pseudo-relevance feedback** Instead of collecting real feedback for a query, we assume that the top  $k$  results are correct, we then use them as a feedback for the Rocchio algorithm.

To motivate this choice we can look at the precision and recall curve, indeed a small set of results usually lead to high precision, thus we can move the system in the direction of the top results.

**Indirect relevance feedback** In this case the feedback is obtained by looking at user behavior, for instance we could count the number of visits to the results to measure the popularity of documents.

## 3.2 Global Methods

**Query expansion** The idea here is to expand the query by adding new information, this could mean adding terminology or shifting the query vector as before.

### 3.2.1 How to add information?

To add information to a query we can perform *syntactic analysis* of the original query, using *dictionaries* to add new terms related to the query, or use the *corpus of reference documents* to add terminology.

**Dictionaries and knowledge bases** The process to enrich a query when having a knowledge base, that could be a dictionary an ontology or whatever, is to:

1. *Lookup* the original query on the external knowledge base
2. Extract *candidate terminology*, so terms that could be added to the query and are somehow related to the one in the original string
3. Use a *word sense disambiguation* tool to split relevant and non relevant terms, and add the relevant ones to the original query



**Wordnet** As introduced in section 1.2 on page 3, Wordnet can be used to extract lemmas but also to find hypernyms and hyponyms of a given word. For instance, if a query presents the word *play*, with the meaning of a *dramatic composition* we could add the latter to the original query.

Another use case could be searching for something like *President Lincoln*. In this case we check for the definition of the two terms, and find out that there's a matching in some possible definitions, indeed the first one could be related to the president of the US, and also the second, thus we can infer that the query is related to US presidency and enrich the query with this information.

**Statistical terminology expansion** Given the set of relevant and non relevant documents  $R$  and  $N$  we can select relevant terminology by using the *pointwise Kullback-Leibler divergence*. We compute the probability of a word to be in  $R$  and in  $N$

$$p(w) = \frac{\text{count}(w, R)}{\sum_{w_i \in R} \text{count}(w_i, R)}, \quad q(w) = \frac{\text{count}(w, N)}{\sum_{w_i \in N} \text{count}(w_i, N)}$$

and then compare them

$$\delta_w(p \parallel q) = p(w) \log \frac{p(w)}{q(w)}$$

if this quantity is close to zero a term is not useful to separate relevant and non relevant documents. On the other end, words with a positive  $\delta$  can be used to expand the query.

### 3.2.2 Spelling correction

There are main phases in spelling correction, the first one is to select the correct versions of a certain query in a set of candidates and then choosing the best one, i.e. the most common or the nearest one.

To detect a *spelling error* we could see few query results, or maybe search for the query tokens in a vocabulary.

**String similarity** A measure that computes the distance of two strings, one could be the *Levenshtein distance*, also known as the edit distance. It can be solved with dynamic programming implementing the following recursive definition:

$$DP[i][j] = \begin{cases} j & \text{if } i = 0 \\ i & \text{if } j = 0 \\ DP[i-1][j-1] & \text{if } s[i] = r[j] \\ \min(DP[i-1][j], DP[i][j-1], DP[i-1][j-1]) + 1 & \text{otherwise} \end{cases}$$

**Phonetic correction** To minimize errors due to phonetic misspelling the idea is to use a *phonetic hash*, so that similar sounding words have the same value. These algorithm are called *Soundex* and are language-dependent.

**K-gram classification and matching** The idea is to consider for each word the subsequences of  $K$  characters.

An example, for the word *PLAY*, is to first add a start and end symbol,  $\#s$  and  $\#e$  respectively and then computing the subsequences, counting how many times they appear in the word.

This creates a vector representation for the words, and to find similar words we can use cosine similarity like we did in section 1 on page 2.

**Sequence-to-sequence learning** This technique use deep learning to match a sequence into another one. The idea is to train a model to map misspelled words into they correct counterpart.