

601.466/666 (Spring 2021)

Assignment 2: Vector Models for Retrieval

Due Date Parts 1 and 2: March 3, 2021, at 5:00 PM EST

Due Date Part 3: March 16, 2021, at 5:00 PM EST

Overview

The goal of this assignment will be to build a vector-based IR engine, similar in underlying design to the Salton/Cornell SMART system. Some of the basic infrastructure for this project has been provided. You will explore the effects of several permutations and implement optional extensions to the model.

The data used for this assignment will be the standard CACM abstract collection, consisting of 3204 computer science journal abstracts from the years 1958-1979. In addition, you will have access to 33 queries for this collection and associated relevance judgements for training and development purposes. Additional queries and relevance judgements have been withheld as unseen test data. Although this collection has several weaknesses, its relatively small size and familiar subject area make it reasonably well suited for rapid experimentation and model development.

Infrastructure

You have been provided with a program `hw2.py` which implements a retrieval engine with several key parts missing. For the data files, the corpus is in `cacm.raw`, the queries are in `query.raw` in the same format, and a list of relevant documents per query is in `query.rels`. There may be some extra files from previous versions of the homework which could be useful.

Corpus processing tools

Several tools for processing the corpus are built into the code for the assignment:

- tokenization: use `word_tokenize` from NLTK
- stemming: use `SnowballStemmer` from NLTK
- stopwords: we provided a list in `common_words`
- document frequency: for computing tf-idf, use the `doc_freq` method in the code

Part 1 - Vector Retrieval Model

You will make several additions/modifications to the program `hw2.py` to implement key components of the vector retrieval model. In particular, you will need to implement the term weighting methods and

similarity measures outlined in Part 2, the precision-recall calculations (methods discussed in class) and functions for system performance feedback. Specifications and suggestions for implementation will be provided in the code.

For your reference:

$$\begin{aligned}
Cosine_Sim(Doc_i, Doc_j) &= \frac{\sum_{t=1}^T (wt_{i,t} \cdot wt_{j,t})}{\sqrt{\sum_{t=1}^T (wt_{i,t})^2 \cdot \sum_{t=1}^T (wt_{j,t})^2}} \\
Dice_Sim(Doc_i, Doc_j) &= \frac{2 [\sum_{t=1}^T (wt_{i,t} \cdot wt_{j,t})]}{\sum_{t=1}^T wt_{i,t} + \sum_{t=1}^T wt_{j,t}} \\
Jaccard_Sim(Doc_i, Doc_j) &= \frac{\sum_{t=1}^T (wt_{i,t} \cdot wt_{j,t})}{\sum_{t=1}^T wt_{i,t} + \sum_{t=1}^T wt_{j,t} - \sum_{t=1}^T (wt_{i,t} \cdot wt_{j,t})} \\
Overlap_Sim(Doc_i, Doc_j) &= \frac{\sum_{t=1}^T (wt_{i,t} \cdot wt_{j,t})}{\min(\sum_{t=1}^T wt_{i,t}, \sum_{t=1}^T wt_{j,t})}
\end{aligned}$$

where T is the total number of terms, Doc_i and Doc_j are two vectors, and $wt_{i,t}$ is the weight of term t in document vector i .

Implementation of term vectors

Although term vectors conceptually have a length equal to the vocabulary size (12,000 terms in this domain), on average only 10-200 of the terms have non-zero values in any given vector. Thus, in practice it is often better to use a dictionary for implementing term vectors. Term weights in a vector can be retrieved efficiently through hashing (e.g. `vec[term]`, with non-present terms returning a weight of 0. The sample code implements cosine similarity using such a dictionary.

Suggestion: the computation of cosine similarity can be made more efficient by precomputing and storing the sum of the squares of the term weights for each vector, as these are constants across vector similarity comparisons.

Implementation of recall/precision calculation

When comparing various permutations of the vector models, it is useful to compare precision at different fixed levels of recall, and also to compute a single measure of mean precision at different levels of recall. The methods used for computing precision and recall on small relevant sets will be discussed in class. For standard comparison of results, we will use the mean precision at 3 fixed levels of recall, averaged over 10 levels of recall, and two normalized measures:

$$\begin{aligned}
Prec_{mean1} &= \frac{Prec_{0.25} + Prec_{0.50} + Prec_{0.75}}{3} \\
Prec_{mean2} &= \frac{1}{10} \sum_{i=1}^{10} Prec_{(Rec=\frac{i}{10})} \\
Recall_{norm} &= 1 - \frac{\sum_{i=1}^{Rel} Rank_i - \sum_{i=1}^{Rel} i}{Rel (N - Rel)}
\end{aligned}$$

$$Prec_{norm} = 1 - \frac{\sum_{i=1}^{Rel} \log Rank_i - \sum_{i=1}^{Rel} \log i}{\log N! / ((N - Rel)! (Rel)!)}$$

where Rel is the total number of relevant documents for the query, N is the total number of documents in the collection, and $Rank_i$ is the position of the i th relevant document in a rank ordered list of all documents sorted by their expected relevance to the query. These last two measures quantify the deviation from the optimal rank ordering where the relevant documents are the top Rel entries in the retrieved list.

Note that the formula for $Prec_{norm}$ (normalized precision) has several large factorials. A reasonable method of handling these is to use the approximation that $\log(n!) = n \log n$ and compute the denominator $\log(N! / ((N - Rel)! (Rel)!))$ as $\log N! - (\log(N - Rel)! + \log(Rel!))$ which can be approximated as $N \log N - (N - Rel) \log(N - Rel) - (Rel) \log(Rel)$.

Because the denominator is a constant across all different permutations of the algorithm (e.g. term weighting schemes, similarity measures, etc.), computing its value exactly is not necessary for a useful comparison of performance across methods.

Part 2 - Experiments in modifying parameters of the vector model

The goal of this part of the assignment will be to achieve both a quantitative measure and an intuitive feel of the effects of changing certain model parameters.

For each permutations of model parameters described below, compute the precision/recall measures described above, averaged over all 33 queries. It is useful to represent the output in tabular form:

Permutation Name	$P_{0.25}$	$P_{0.50}$	$P_{0.75}$	$P_{1.00}$	P_{mean1}	P_{mean2}	P_{norm}	R_{norm}
Raw TF weighting								
TF IDF weighting								
Boolean weighting								
Cosine similarity								
Dice similarity								
...								

The starter code already prints out a table in this format.

1. Term weighting permutations:

- (a) Raw TF weighting ($wt_{t,d}$ = raw frequency of term t in document d).
- (b) * TF-IDF weighting ($wt_{t,d} = TF_{t,d} \cdot \log(\frac{N}{DF_t})$)
- (c) Boolean weighting ($wt_{t,d} = 1$ if term t is present in doc d , 0 if absent)

2. Similarity measures:

- (a) * Cosine similarity
- (b) Dice, Jaccard, and Overlap similarity

3. Stemming

- (a) Use raw, unstemmed tokens (all converted to lower case)
- (b) * Use tokens stemmed by the Porter stemmer

4. Stopwords

- (a) * Exclude stopwords from term vectors
- (b) Include all tokens, including punctuation

5. Region weighting

- (a) Weight titles, keywords, author list and abstract words equally
- (b) * Use relative weights of titles=3x, keywords=4x, author list=3x, abstract=1x.
- (c) Use relative weights of titles=1x, keywords=1x, author list=1x, abstract=4x.

The provided starter code iterates through each of these permutations, printing out the precision/recall measures in a tabular format. Please produce a file containing a table of results with the following command:

```
python hw2.py > output.tsv
```

Writeup

Produce a writeup detailing your thoughts and observations for these permutations. Please comment on permutations that performed particularly well or poorly, and why you think they did so. You may also comment on when you might pick one permutation over another.

In addition, for a qualitative feel of parameter effects, compute and submit the following output for combinations 3a and 3b above. You are strongly encouraged to examine results for other permutations and other queries on your own.

1. List the top 20 retrieved documents for Queries 6, 9 and 22 by their number, title and similarity measure, with the “relevant” documents starred.
2. For the top 10 retrieved documents, show the terms on which the retrieval was based (those with non-zero weights for both query and retrieved document) along with these weights.

It is interesting to see how little information actually contributes to identifying document relevance.

3. List the top 20 documents that are most similar to Documents 239, 1236 and 2740, giving number, title and similarity measure.

Part 3 - Extensions to the retrieval model

Implement extensions to the retrieval model selected from the set below, preferably original ones that improve performance over that observed in the given permutations above. Enter the results in precision/recall table above, labelled as permutations 6 (and 7 if necessary). Since different extensions have different complexity and expected time cost, rough estimates for this complexity are given below. Your total complexity score must be at least 3.

Note: If your total complexity score is 4 or greater you may significantly shorten Part 2 by only providing evaluation data for 2 permutations (the default (starred) permutation, and a second contrastive permutation 1a,2b,3a,4b,5c) plus results for any option you implement below, and you can skip the implementation of P_{norm} and R_{norm} .

If an extension does not have an obvious empirical evaluation, you do not need to provide one.

1. Implement a simple relevance feedback method, computing a new query centroid vector from a weighted linear combination of the original query, the relevant-labelled documents and the (negatively weighted) irrelevant-labelled documents. Note that the person giving the relevance feedback (you) may have different standard of relevance than the original query formulator.
You can do an automatic (albeit biased) test of your performance by using the true relevance judgements (given in `query.rels`) on the top 20 documents initially ranked by your system as if these judgements were coming from the user. This feedback can then be used as above to re-rerank all the data in the 2nd round. **[complexity: 2]**
2. Create a new corpus of documents from any document set of your choice. This can be web pages, news stories, email, files in your home directory, literally any collection of things that you may wish to search. The collection should have at least 100 documents, and preferably more than 200. You should provide at least 10 plausible queries on this document set and provide relevance judgements for these queries. **[complexity: 2.5]**
3. Add 3 meaningfully improved term weighting strategies of your choice, possibly including a more refined weighting by region of document and negative term weighting for regions of queries specifying subjects/sub-areas that should *not* be included, such as in query 18. **[complexity: 1]**
4. Augment the term set to include all *bigrams* in the document/query that do not contain stopwords. For example “*window managers and command interpreters*” would add the bigrams *window+managers* and *command+interpreters* to the term set, but not *managers+and*. This will help make searches sensitive to word order rather than unordered, bag-of-words co-occurrence.
The associative-array representation of vectors should make this easy to implement. You may use `make_hist.pl` to compute overall document frequencies by giving it a stream of bigrams (easy to generate). **[complexity: 1.5]**
5. Redefine the term set to include all 5-character n-grams in the text (including spaces (as “_”) and hyphens, excluding other punctuation), and including padded spaces at the beginning and end of queries. For example, a query (or document) with the string *window managers* would contain the terms “_wind”, “windo”, “indow”, “ndow_”, “dow_m”, “ow_man”, etc. This technique tends to be robust in the face of OCR errors and other noise, and effectively captures some multi-word terms. You may remove stopwords first at your discretion, but please state if you do so. **[complexity: 1.5]**
6. Implement query expansion through the use of thesaurus classes. **[complexity: 1.5]**
7. Cluster the top k returned documents (or those above similarity s) using the basic greedy Salton clustering method described in class. List documents in each cluster separately, with a line separating the clusters and ordered by similarity within clusters. The number of clusters you choose can be based on any reasonable criteria. Label each cluster with the most salient terms found in each cluster (for example, high frequency and high TF-IDF terms in the cluster). This would be useful when a query such as “windows” returns clearly different partitions of documents on (Microsoft) windows and (glass) windows. **[complexity: 2]**
8. Deploy and compare 2 word embedding methods (e.g. using GloVe) via either term clustering or query expansion. **[complexity: 2.5]**
9. Reduce the dimensionality of the term vectors using SVD. **[complexity: 3]**
10. Modify your program to accept queries directly from the keyboard, rather than requiring pre-computation of a query corpus. **[complexity: 2]**

11. Create a basic web interface to your search engine. You should read in a query from an HTML form, call a modified vector.pl program as a CGI script, and generate output as HTML, with simple HTML links to the documents (you can use "#" offsets into the full abstract file, no need to split into separate document files). It is not necessary to implement persistent state for your program or any special loading efficiency, although implementing some of the large data structures as rapid-load DBM files would help speed performance. **[complexity: 2]**
12. Create the web interface as described above, but implement your search engine as a persistent server that does not need to be invoked for each new query. You should support multiple socket connections and a mechanism for session tracing. Ideally you should also support feedback buttons for "more like this" relevance feedback. This should effectively act as a full web-based search engine. This option should only be chosen for those who have the background to implement such socket connections, etc. on their own. **[complexity: 2 or 2.5 (with relevance feedback interface)]**

When appropriate, your two extensions will also be evaluated using a second set of held-out queries. Thus methods should not be tailored specifically to idiosyncrasies in the given query set. The test queries and relevance judgments are from the same source as the training material, however, and have similar properties.

Programming

Similar to Homework 1, we must be able to run your code using the following command:

```
python hw2.py
```

For Parts 1 and 2, you may not use any external libraries or tools that do the homework for you. For example, `numpy.linalg.norm` is fine, but `sklearn.metrics.pairwise.cosine_similarity` is not. This restriction is lifted for Part 3. For example, you may use `numpy.linalg.svd` to perform SVD if you wish.

Evaluation

Submissions will be evaluated as follows:

Part 1: Implementation of the vector-based retrieval engine	40%
Part 2: Exploration of permutations and writeup	15%
Part 3: Quality and creativity of your extensions to the model	40%
Part 3: Performance of your extensions on held-out test data	5%

Submission

Homework submission will be via Gradescope. Please follow the same instructions as for Homework 1, except replace hw1 with hw2.