

INFORMATION RETRIEVAL

Flavio Forenza

flavio.forenza@studenti.unimi.it



Department of Computer Science,
University of Milan, Italy

September 19, 2021

INTRODUCTION

Nowadays, searching the web for information appears to be one of the simplest operations to perform. The difficulty perceived by the user in formulating a query has been gradually reduced by techniques capable of guiding his writing towards a correct generation of a query. These techniques allow to improve the performance of information search systems.



GOAL

The purpose of this project is to be able to experiment with the use of one of the most famous techniques, already present at the state of the art, able to "assist" the user in formulating a correct query: **Language Modelling**. Query expansion can be done using this concept to return a corpus of relevant documents.

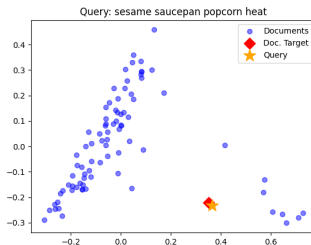
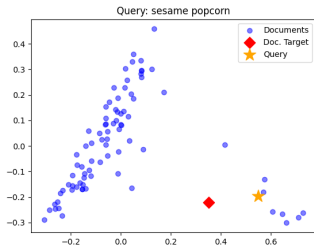


Figure. Distance between initial query and the target document in a PCA space.

Figure. Distance between the expanded query and the target document in a PCA space.



DATASET DESCRIPTION

The dataset used for the experiments is the famous *Recipes1M+*¹, a collection created by MIT, consisting of more than one million culinary recipes. Of all these recipes, only a subset of 51235 documents of it was used due to their informative content which best fits the purpose of this study. The information about the line distributions for each recipe indicates that the instruction field contains a higher number than the information contained in the ingredients field.



Figure. Distributions of lines per ingredients and instructions.

¹J.Marín, A.Biswas, F.Ofli, N.Hynes, A.Salvador, Y.Aytar, I.Weber and A.Torralba, "Recipe1M+: A Dataset for Learning Cross-Modal Embeddings for Cooking Recipes and Food Images", IEEE Trans. Pattern Anal. Mach. Intell., 2019



RANKING GENERATION

The first step is based on choosing a random query*, which resembles the title of one of the existing recipes (the *target document*). Subsequently, thanks to the combination of the tf-idf method and the cosine similarity metric, it is possible to generate the first ranking of documents ordered according to relevance with the query. The threshold chosen, for the selection of the most relevant documents, will correspond to the weight assigned by the tf-idf to the target document.

tf-idf

$$TfIdf(q_t, d) = tf_{q_t, d} \log \frac{N}{df_{q_t}}$$

Cosine similarity

$$cosine(q, d) = \frac{\sum_{i=1}^N qd_i}{\sqrt{\sum_{i=1}^N q^2} \sqrt{\sum_{i=1}^N d_i^2}}$$

* All documents (instructions of each recipe) and queries (recipe title) will undergo tokenization and normalization processes (removing stopwords / punctuation and applying lemmatization) before generating the ranking.



RANKING EVALUATION

The evaluation of the ranking of relevant documents was carried out considering as the entity, of each single document, its category of belonging² (dessert, salad, beverage, etc.). The search for each category takes place in two methods, with two different libraries:

- *Scrape Schema Recipe*
- *USDA* (United States Department of Agriculture)

Queries	Scrape Schema Recipe			USDA	Mixed (Scrape+USDA)
	Overestimate	Underestimate	Discarded		
I	0.7520	0.3734	0.5958	0.9817	0.9947
II	0.9309	0.6780	0.9054	1.0	1.0
III	0.7458	0.2746	0.5411	0.9797	0.9982
IV	0.8939	0.5320	0.8433	0.9612	0.9870
V	0.8335	0.5033	0.7544	0.9940	0.9940
Average	0.8312	0.4722	0.7280	0.9833	0.99478

Table: Average Precision on each query for each method.

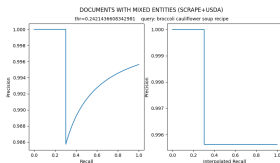


Figure. Performance in terms of Precision, Recall and Interpolated Recall.

²Md R.Parvez et al. "Building Language Models for Text with Named Entities". In: (2018), pp. 373–2383.



LANGUAGE MODELS

For each document, inside the ranking, a language model is built. Each language model is initially composed of a sequence of terms equal to two (*bi-grams*) and their occurrence. Subsequently it will be possible to calculate the *sentence probability* (1) between the query and each single language model.

Sentence probability (Bi-gram) and MLE

$$\begin{aligned} P(q|d) &\approx P(q|M_d) \\ &\approx \prod_i^n P(w_i|w_{i-1}) \\ &\approx \prod_i^n \frac{\text{count}(w_i, w_{i-1})}{\sum_{j=1}^n \text{count}(w_j, w_{i-1})} \\ &= \prod_i^n \frac{\text{count}(w_i, w_{i-1})}{\text{count}(w_{i-1})} \end{aligned} \tag{1}$$



SMOOTHING METHODS

To avoid having a probability equal to zero, two smoothing techniques are calculated:

Laplace Smoothing

$$P(w_i|w_{i-1}) = \frac{\text{count}(w_{i-1})+1}{\text{count}(w_{i-1})+|V|}$$

where:

- $\sum_i \lambda_i = 1$
- M_d : LM of the single document;
- M_c : LM of the entire collection of documents;
- $|V|$: # unique words within the corpus of documents.

Linear Interpolation Smoothing (Bi-grams)

$$P(w_i|w_{i-1}) = \lambda P(q|M_d) + (1 - \lambda)P(q|M_c)$$

⋮

Linear Interpolation Smoothing (Zero-grams)³

$$P(w_i) = \lambda \frac{1}{|V|} + (1 - \lambda)P(w_i)$$

³ A.Gutnik, "Log-Linear Interpolation of Language Models", 2000



CORE

Algorithm 1: Best Ranking

```

input      : q
output    : final_ranking
parameter : LMtd

for s ← 2 to 10 do
    LMdocs ← ComputeLMdocs(s, query);
    LMcoll ← ComputeLMcoll(s, query);
    rankLpl ← Lp1Smoot (q, LMdocs);
    perplexityLpl ← perplexityLpl (q, LMtd);
    lamb1, lamb2 ← 0;
    for i ← 0.1 to 1.0 do
        rankLinInt ← LinInt (q, LMdocs, LMcoll);
        perplexityInt ← perplexityInt (q, LMtd);
        lamb1 ← 0.1 + i;
        lamb2 ← 0.9 - i;
        if lamb2==0 then
            break;
        end
    end
end

end
.
.
.
Search for the minimum index of the target document;
final_ranking ← min(rankLpl[i],rankLinInt[i]);
return final_ranking
    
```

where:

- *q*: query
- *s*: skip-gram
- *final_ranking*: ranking containing the target document in the position closest to the first
- *LMtd*: Language model of target document
- *LMdocs*: List of the language models of each document
- *LMcoll*: language model of the entire collection of documents
- *lamb1*: λ_1
- *lamb2*: λ_2



TERM-TERM MATRIX

As explained in⁴, two or more words are considered synonymous, or semantically important, if their *vectors*, represented in a multidimensional space, have a high *cosine similarity*. To find out, we need to build a matrix of terms, called *term-term matrix*, where the *co-occurrences* between the terms present in the relevant documents and in the query will be reported. We can exploit the *Language Model* of the entire corpus, of relevant documents, to derive co-occurrences.

	variation	broccoli	cauliflower	toss	substitute	...
variation	0	1	0	0	0	...
broccoli	0	0	5	0	0	...
cauliflower	0	0	0	1	0	...
toss	0	0	0	0	1	...
substitute	0	0	1	0	0	...
...

Table: Co-Occurrence of words in Term-Term Matrix.

⁴D.Jurafsky and J.H. "Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition". In: vol. 3. 2008.



POSITIVE POINTWISE MUTUAL INFORMATIONS (PPMI)

Pointwise Mutual Information (PMI) was used to evaluate the relationship between two words.

PMI

*"PMI draws on the intuition that the best way to weigh the association between two words is to ask how much more the two words co-occur in our corpus than we would have a priori expected them to appear by chance."*⁴

In this case, its "Positive" version (PPMI) was used to be able to remove negative values.

PMI

$$PMI(x, y) = \log_2 \frac{P(x, y)}{P(x)P(y)}$$

PPMI

$$PPMI(x, y) = \max(\log_2 \frac{P(x, y)}{P(x)P(y)}, 0)$$



SINGULAR VALUE DECOMPOSITION (SVD)

The term-term matrix presents two types of problems that affect the performance and the calculation of cosine similarity: **high dimensionality** and **sparsity**. In order to reduce the dimensionality, the *Singular Value Decomposition (SVD)* technique was applied.

Where:

$$A = USV^T$$

$$\mathcal{D} = U * S$$

$$\mathcal{T} = S * V^T$$

- A : a $t \times d$ Term-Term Matrix decomposed into 3 sub-matrices.
- U : a $t \times m$ matrix
- S : a $m \times n$ matrix
- V^T : a $m \times d$ matrix
- \mathcal{D} : a $t \times m$ matrix containing all the terms of the relevant documents
- \mathcal{T} : a $m \times d$ matrix containing the query terms useful for calculating the cosine similarity with the terms present in \mathcal{D}



QUERY EXPANSION

After calculating the cosine similarity between the word vectors of the relevant documents and the word vectors of the query, a maximum of ten most similar words are selected for each single word in the query. Starting from this assumption, by making a combination between the pairs of words (q_t, d_n) , the maximum number of new queries that can be generated is:

$$\#Queries = \prod_i^q (\#sim_words_i)$$

Where q represents the number of tokens of the query and $\#sim_words_i$ represents the number of tokens most similar to each single token i of the query. Obviously some checks are made (eg word repetition), which could invalidate the generation of the query. **All new queries will be processed by the algorithm (Core) in order to determine the best ranking of relevant documents, which has the target document in the first positions.**



PERPLEXITY

The perplexity is calculated between the initial query/new_queries and the language models of the target document. The peculiarity highlighted by the system is that as the number of skip-grams increases, as the lambda parameters of Linear Interpolation Smooth- ing vary, the perplexity gradually decreases.

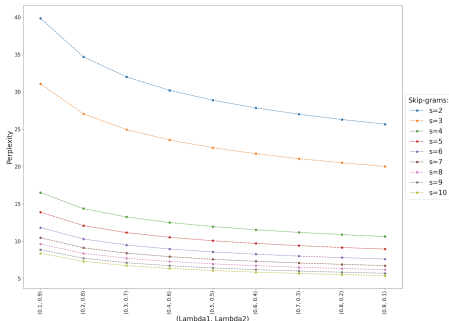


Figure. Variation of the perplexity value based on the change of the step s and of the λ_1 and λ_2 values of interpolated smoothing.



SYSTEM EVALUATION WITH DIFFERENT PARSERS

To evaluate the performance of the system, four different parsers were used: **Keras**, **Nltk**, **Spacy** and **Gensim**. The search results of the target document are compared with those obtained by the classic method *tf-idf*. Each system has its own benefit but all of them manage, among the first 100 queries generated, to place the target document in the first position of the ranking.

Method	I		II		III		IV		V	
	T	P	T	P	T	P	T	P	T	P
Keras	228	0	623	0	82	0	126	0	51	0
Nltk	228	0	623	0	82	0	126	0	51	0
Spacy	221	0	647	0	81	0	141	0	72	0
Gensim	240	0	652	37	84	0	113	0	37	0

Table: Target document position in the ranking of relevant documents. (T: *tf-idf*, P: *Proposed*)

Method	#Queries	Time
Keras	14.900	Low
Nltk	14.900	Low
Spacy	19.010	<i>High</i>
Gensim	22.090	<i>Medium</i>

Table: Performance in terms of number of queries (#Queries) and computation time.



PERPLEXITY AND DISTANCE EVALUATION WITH DIFFERENT PARSERS

Comparing all the systems, **Gensim** manages to have, in most cases, a *low perplexity*, correlated to a *low Euclidean distance*, between a new generated query and the target document. These results make *Gensim* the best performing system.

Queries	Keras		Nltk		Spacy		Gensim	
	D	P	D	P	D	P	D	P
I	0.14	16.41	0.14	16.59	0.12	34.89	0.14	90.17
II	228	17.91	0.11	22.37	0.2	11.57	0.15	53.44
III	0.02	5.41	0.02	5.41	0.07	5.15	0.01	5.03
IV	0.05	12.39	0.05	12.45	0.09	12.44	0.01	10.83
V	0.31	23.04	0.31	23.22	0.4	42.25	0.27	10.36
VI	0.22	6.39	0.22	6.39	0.24	3.93	0.21	4.02
VII	0.12	1429.33	0.12	1521.33	0.24	6.03	0.1	3.43
VIII	0.39	13.8	0.39	13.81	0.46	13.65	0.27	8.11
IX	0.56	5.3	0.56	5.3	0.6	4.7	0.41	3.9
X	0	4.29	0	4.29	0.06	5.12	0	4.28

Table: Euclidean Distance (D) between a query and the target document at a specific perplexity (P).



CONCLUSIONS

Coming to the conclusions, it can be said that the generation of new queries occurs correctly. However, there are some points that need to be improved, such as determining a method that sets an ad **hoc threshold** useful for generating each ranking of relevant documents. Another improvement concerns the reduction of the **sparsity** of the term-term matrix. This problem can be solved by applying techniques such as the *stochastic gradient descent*. As for future work, however, it would be interesting to develop a system of **hybrid parsers**, fast and with low perplexity. In addition to this development, it will be interesting to adopt a new method of query expansion with terms present in a vocabulary such as **Wordnet**.



THANKS FOR LISTENING



▸ Link Project

▸ Link Paper

▸ Link Presentation

