

INFORMATION RETRIEVAL

Flavio Forenza

flavio.forenza@studenti.unimi.it



Department of Computer Science,
University of Milan, Italy

September 11, 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.



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. 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

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



RANKING EVALUATION

The evaluation of the ranking of relevant documents was carried out considering as the entity, of each single document, its category of belonging¹ (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.

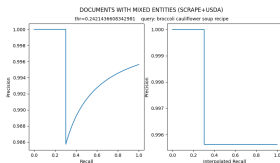


Figure. Performance in terms of Precision, Recall and Interpolated Recall.

¹Md R.Parvez et al. "Building Language Models for Text with Named Entities". In: (2018), pp. 373–2383.



LANGUAGE MODELS

For each document, present in the ranking, having a threshold greater than or equal to the weight assigned to the target document, a language models will be constructed consisting of a sequence of *bi-grams*, with an initial *skip-grams* equal to two, with the relative count of every occurrence.

Sentence probability (Bi-gram) and MLE

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



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



TERM-TERM MATRIX



POSITIVE POINTWISE MUTUAL INFORMATIONS (PPMI)



SINGULAR VALUE DECOMPOSITION (SVD)



QUERY EXPANSION



PERPLEXITY



SYSTEM EVALUATION WITH DIFFERENT PARSERS



CONCLUSIONS

