

# Cluster Rank Demo Harness - Unfortunate Series of Events

Philip Robinson

*Oregon Health Sciences University*

June 26, 2017

## Abstract

It is often the case that initial query compositions result in frequent restarts as the user negotiates with their retrieval system. This is likely a product of unfortunate query formulations or choice of ranking algorithm. Our proposed retrieval system encourages diversity in displayed documents by introducing an unsupervised clustering step before displaying results. The clusters are then presented to the user with their documents ranked independent of each group. We do this by clearly separating the retrieval process into the three steps **relevance, clustering, and ranking**, then allow the user to recurse this process on a cluster (rather than restarting their query). Additionally, we propose a simple method to compare results against varying quality **tfidf** queries. Our final product is a demo harness that abstracts these steps, so that others may easily produce and reproduce prototypes against their own corpora.

## 1 Introduction

Information retrieval systems have long suffered from non-informative query formulations by their users. Many systems employ techniques such as query expansion, domain ontologies, advanced search parameters, to address such difficulties. Unfortunately, most retrieval systems presume user provided queries to be informative, and select to return the most query-relevant documents. We can interpret a query-relevance ranking on documents retrieved by a non-informative query as being an over-fit ranking (to the query). Unfortunately, this over-fitting can happen regardless of query quality. In cases where all documents retrieved are nearly identical.

Either in the case of non-informative queries or of monolithic document listings, users are often tasked with query reformulation. Reformulation may be an insufficient mechanism, for users not appropriately familiar with their target domain. Additionally, as a side effect of this workflow, query reformulation can make it difficult to assess the general effectiveness of a retrieval system [1].

To prevent over-fit rankings some retrieval systems have proposed document diversification strategies [2]. This can be accomplished by introducing similarity or taxonomic clustering of documents

either prior to or after query submission. Additionally, to improve query formulation many search engines provide similar search terms along side retrieved documents. These aid terms are usually assigned by experts to similar vocabulary or to retrieved documents, and can be expensive to curate.

## 1.1 Proposal

We propose to develop a harness for prototyping information retrieval systems against corpora where clustering and ranking do not inform each other. In our system, we propose separating the query processing pipeline into quarantined steps. We first identify documents by relevance, then perform unsupervised clustering of relevant documents, finally rank each cluster by the query provided. We then allow the user to zoom in on a cluster. During this process, every cluster is also tagged with terms automatically, by providing the grouping’s most impactful `tfidf` words. These terms don’t impact the ranking, but act to provide the user with additional information in selecting a term.

We hypothesis that a document clustering will allow users to eliminate poor document groupings. Additionally automatic tagging of clusters with relevance terms should allow users to navigate a listing retrieved from non-informative queries. To form queries, we sample words from randomly selected documents’ `tfidf` distributions and capture the resulting ranking  $R_{\text{tfidf}}$ . For a system with  $C$ -way clustering to be considered permanent, the same query should yield our target document within  $S$  steps.

$$S < \log_C(R_{\text{tfidf}})$$

As this is a process greatly empowered by composable components, we also separate out these concerns in code, to hopefully decrease testing time against multiple corpora, ranking, and clustering configurations.

## 2 Implementation Details

We provide a retrieval systems harness, developed in `python`<sup>1</sup>. The harness is designed to be used from the command line, but hosts results in a manner accessible to a flask server (that is also provided).

In order to reduce complexity<sup>2</sup>, `numpy` arrays were used to track, sort, and split document collections. Interfaces that required of designers must yield index arrays. Unfortunately, (presently) this limits corpus scaling to what can be handled by `numpy`. As an additional restriction, corpora are required to be transformed to `numpy` compliant structures.<sup>3</sup>

---

<sup>1</sup><https://github.com/probinso/IR-cluster-rank-demo>

<sup>2</sup>and scratch the author’s wish to use array languages

<sup>3</sup>examples found in code

Provided is the `olib.Handler` class that orchestrates the cluster, ranking, steps as a coroutines. An `olib.Organizer` is passed as a constructor to the `olib.Handler`. `olib.Organizers` which are `numpy.ndarrays` extended with interfaces for `.cluster(vocabulary, n_clusters)` and `.rank(*rankargs)`, which raise `NotImplemented` unless provided by designer. Designers implement their own `olib.Organizers` and overload these methods.

For our test cases, we have `kmidf.CosKMIDFROrganizer` which is a `tfidf` ranker coupled with k-means over cosine distance, and `ldaidf.LDAIDFROrganizer` which is a `tfidf` ranker coupled with Latent Dirichlet Allocation max clustering. Classes pretended with `Augmented` have been altered for evaluation purposes.

Modifications to a handler allow us an oracle into the performance of a system. `olib.AugmentedHandler` tags a cluster's meta-data with ownership of a target document, and it's rank location. This allows us to use traverse a guided path through our documents. We still require that discovery of a document is measured by it's being placed in the displayed range (top 3 by default) for a cluster.

```
1 class IDFRankOrganizer(Organizer):
2     def rank(self, queryvec):
3         dist = np.sum(self[:, queryvec],axis=1)
4         idx_dist = sorted(enumerate(dist),
5                           key=itemgetter(1),
6                           reverse=True)
7         idx = np.array([key for key, value
8                         in idx_dist])
9         return idx
```

```
1 class CosKMOrganizer(Organizer):
2     def cluster(self, terms, n_clusters=3):
3         dist = 1 - cosine_similarity(self)
4         alg = KMeans(n_clusters=n_clusters)
5
6         # Fit cosine distance kmeans
7         results = alg.fit(dist)
8         labels = results.labels_
9         complete = dict()
10        ulabels = np.unique(labels)
11        for l in ulabels:
12            keys = np.array([i for i, b in
13                            enumerate(labels==l)
14                            if b])
15            complete[l] = keys
16
17        ... # populate meta
18
19        return complete, meta
```

```
1 class LDAOrganizer(Organizer):
2     def cluster(self, terms, n_clusters=3):
3         lda = LatentDirichletAllocation(
4             n_topics=n_clusters)
5         lda.fit(self)
6
7         topics = lda.transform(self)
8
9         # argmax is a bad assignment choice
10        assign = topics.argmax(axis=1)
11        labels = np.unique(assign)
12
13        complete = dict()
14
15        for l in labels:
16            complete[l] = np.where(assign==l)[0]
17
18        ... # populate meta
19
20        return complete, meta
```

```
1 class CosKMIDFROrganizer(
2     CosKMOrganizer, IDFRankOrganizer):
3     pass
4
5 class LDAIDFROrganizer(
6     LDAOrganizer, IDFRankOrganizer):
7     pass
```

## 3 Evaluation

### 3.1 Approach

In order to provide repeatable results without depending on expertly relevance and ranked corpora, we base our comparisons against traditional `tfidf`. Relevance is determined as any document

containing any of the query terms, and its rank is defined as  $R_{\text{tfidf}}$  denoting its place in the listing, from **tfidf** score against these documents. As mentioned above, for a system with  $C$ -way clustering to be considered proffered, an identical query should yield our target document within  $S$  steps.

In this case, the clusters contain unique document lists<sup>4</sup>. This means that success is measured by satisfying bounds for  $S < \log_C(R_{\text{tfidf}})$ . Using the augmented structures, described above, clusters are labeled to contain the target document. This however would not be the case in a real environment. It is also important to note that some documents contained the only instance of a word in the corpus. To prevent sampling these words, we eliminated words with **tfidf** scores  $\geq 1$  from our sample space.

We generated tests against two different systems, the `kmidf.CosKMIDFROrganizer` and the `kmidf.LDAIDFROrganizer`, with a collection of tenthousand medical abstracts.

The unit retrieval was a structure of cluster labels and corresponding document labels with meta-data. The image below is an example of a first response on `kmidf.CosKMIDFROrganizer`.

```

1 { 0: {
2   'documents': [
3     'Molecular Mechanisms Governing Cooperating Motors',
4     'A wrinkle in time: The perception of expressions in the elderly',
5     'Computer simulations of populations of mammalian motor units'],
6   'meta': [
7     'epitheli', 'knockout',
8     'overexpress', 'mutant',
9     'homeostasi', 'downstream',
10    'extracellular', 'kinas',
11    'induct', 'stem',
12    'HERE!!! :27']],
13 1: {
14   'documents': [
15     "The Cyclical Lower-extremity Exercise for Parkinson's Trial",
16     'OLDER BREAST CANCER PATIENTS: RISK FOR COGNITIVE DECLINE',
17     'Statewide Intervention to Reduce Early Mortality in Persons with Mental Illness'],
18   'meta': [
19     'extracellular', 'intracellular',
20     'knockout', 'mutant',
21     'overexpress', 'transgen',
22     'apoptosi', 'transcript',
23     'mammalian', 'ligand']],
24 2: {
25   'documents': [
26     'Optic Nerve Head SDOCT Imaging in Glaucoma',
27     'Multimodal Connectivity for Surgical Mapping',
28     'Molecular Imaging Directed, 3D Ultrasound-guided, Biopsy System'],
29   'meta': [
30     'apoptosi', 'block', 'cycl',
31     'diabet', 'dna', 'elucid',
32     'endotheli', 'homeostasi',
33     'host', 'immun']]
34 >> 0

```

---

<sup>4</sup>It is not a restriction of the system that clusters need be unique

### 3.2 Results and Discovered Limitations

After selecting many documents, from the corpus, at varying query strengths we were rarely able to satisfy the success criteria. This may have been due to many factors, some of which arose as a serious consequence of pure segregation of clustering from ranking. Ranking that isn't informed by clustering would retain the same top documents, unless you were (randomly) lucky enough split along the top ranked documents. This behavior resulted in stagnating documents at the top of reported clusters, suppressing the target document from being selected.

This document suppression was especially bad in k-means. This was likely because the the scoring metric for document similarity was in the same space as the ranking. Consequently, if a target document was close to a dense centroid, then it would likely never hit the visible documents list; as its competing documents would likely continue to score high from that region. This was a poor outcome.

Additionally, the soft clustering (found in Latent Dirichlet Allocation) introduced problems with reasonable group assignment. As groups became smaller, topics became less informative and would often break down. Complicated group assignment algorithms (informed by LDA) were used, but resulted in meaningless splits nearing the end of a query session.

The most tested document, `id(9871)`<sup>5</sup>, was selected due to the evaluator's ability to understand it's content. The hope was to use a cluster's meta-data, in the form of cluster informed keywords, to guide sub cluster selection. Unfortunately, due to document suppression, dense documents clusters what overlapped with the query results also overloaded the meta search terms. Finally, `tfidf` distributions weren't updated between clustering steps, which resulted in stabilization of terms in groups, (eventually leading to many shared terms across groups). The evaluator was unable to appraise cluster selection without the oracle by use the assisted query terms; likely due to the costs of document suppression. Similar results were found with other documents.

doc.id	Query	$R_{\text{tfidf}}$	LDA ( $S, R_S$ )	Kmeans ( $S, R_S$ )
9871	motor younger medic vivo	52	0:15, 1:12, 2:12, 3:11	0:31, 1:27, 2:19, 3:10
			0:31, 1:30, 2:30, 3:30	0:31, 1:27, 2:17, 3:5
			0:31, 1:27, 2:27, 3:26	0:31, 1:27, 2:17, 3:5
9871	motor age medic vivo	25	0:16, 1:14, 2:12	0:12, 1:1
			0:9, 1:9, 2:9	0:12, 1:1
			0:6, 2:6, 3:6	0:12, 1:1

Interestingly, since both k-means over cos-distance and LDA provided very stable results on this corpus. It has become apparent, that stable clustering may significantly reduce the effects that ranking may have when selecting a cluster to recurse on. Further work may lead to better strategies for perturbing these clusters, so that ranking can have more interesting effects.

---

<sup>5</sup>Molecular Aging of the Human Brain: Genetic Modulation and Functional Outcomes

### 3.3 Limitations

In order to support more sophisticated document clustering and ranking (like those using word2vec) we will need to loosen our restriction to `numpy` data structures. Additionally, significant speedups could be found in the document relevance lookup, resulting in faster testing of configurations to corpora.

Although this was proposed as an evaluation harness, several design decisions weren't made with evaluation in mind. The `AugmentedHandler` behavior should be moved to first priority in the class structure, as to provide sane defaults for potential users of the system.

## 4 Next Steps

Firstly, we would like to address any concerns raised from the limitations section. We are also interested in providing a full restful API so those who find these systems useful may easily host their corpora without much engineering effort.

We are also interested in extending user search to include selecting multiple clusters and inject new query terms without restarting their search.

## 5 Conclusions

This project will likely die in its current incarnation. The results weren't very promising, and the discovered fixes violate the original premis. We are still interested in starting clean with a retrieval system prototyping harness, but likely starting fresh with our newly gained knowlege.

## References

- [1] Jeff Huang and Efthimis N. Efthimiadis. Analyzing and evaluating query reformulation strategies in web search logs. In *Proceeding of the 18th ACM conference on Information and knowledge management - CIKM '09*. ACM Press, 2009.
- [2] Yadong Zhu, Yanyan Lan, Jiafeng Guo, Xueqi Cheng, and Shuzi Niu. Learning for search result diversification. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval - SIGIR '14*. ACM Press, 2014.