# Cluster Rank Demo Harness

Philip Robinson

*Oregon Health Sciences University*

June 24, 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 seperating 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-relevent documents. We can interpret a query-relevance ranking on documents retrieved by a non-informative query as being an overfit ranking (to the query). Unfortunately, this overfitting can happen regardless of query quality. In cases where all documents retrieved are nearly identical.

Either in the case of non-informated queries or of monolithic document responses, users are often tasked with query reformulation. Reformulation may be an inssuficent mechanism,

for users not appropriately familure 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 overfit rankings some retrieval systems have proposed document diversification strategies[2]. This can easily be done by similarity or ontological 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[3].

# 2 Implementation Details

# 3 Evaluation Approach

# 4 Expirimental Results

# 5 Limitations

# 6 References

---

[1]reference
[2]reference
[3]reference