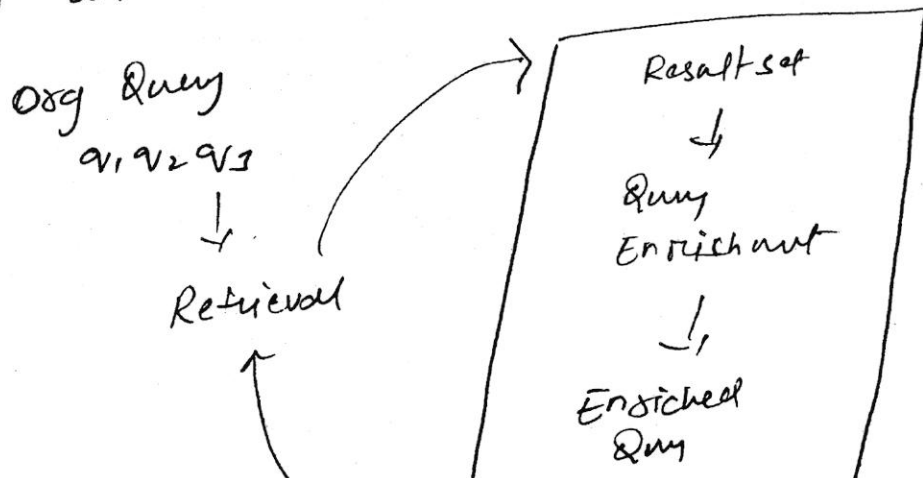


Retrieval Utilities:

Many different utilities improve the results of a retrieval strategy. Most utilities add or remove terms from initial query in an attempt to refine the query.

Relevance feedback: The top doc found by an initial query are identified as relevant. These doc are then examined. They might be deemed relevant either by manual intervention or by an assumption that the top n doc are relevant.

The basic idea is to implement retrieval in multiple passes. The user refines the query in each pass based on the results of earlier query. The user indicates which of the doc presented in response to an initial query are relevant. A new term are added to the query based on this selection. Additionally existing terms can be reweighted based on user feedback.



An alternative approach is to avoid asking the user anything & to simply choose the top ranked doc as relevant. Using either manual or automatic, the initial query is modified & the new query is re-eval.

Ex find info surrounding the various conspiracy theories about the assassination of John F. Kennedy.
has both useful keywords & noise.

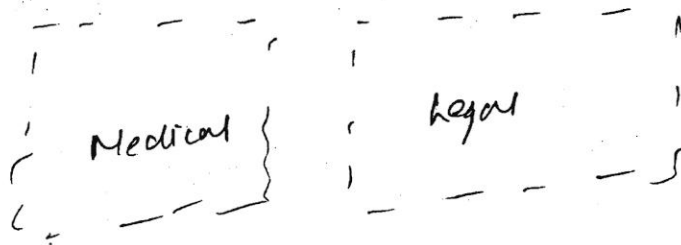
keywords are Assassinate John
various, info are NOT stop words.

frequently used words a, an, and, the are ignored

the idea is to use all the terms in the query & ask the user if top ranked doc are relevant.

Assume a highly ranked document contains the term Oswald. It is reasonable to expect that adding the term Oswald to the query will improve both precision & recall. And if it may contain many occurrences of assassination, the not used in initial query for this term should be increased.

Clustering: doc clustering attempts to group doc 2 (2) by content to reduce the search space required to respond to a query. For ex a documents collection that contains both medical & legal doc might be clustered such that all med doc are placed into one cluster & all legal doc are assigned to legal cluster. A query over legal material might be directed to the legal doc cluster.



Several clustering algorithms have been proposed. In many cases the evaluation has been challenging because it is difficult to automatically point a query to a doc cluster.

① Result Set Clustering:

Clustering was used as a utility to assist relevance feedback. In those cases on top the results of a query were clustered & in relevance feedback process, by day new terms from large clusters were selected.

Recently, web search results were clustered based on significant phrases in the result set. First documents in a result set are parsed & two term phrases are identified. Characteristics about these phrases are then used as input

algorithm. Fed

are used in this work. A cluster is initially identified as the set of documents that contain one of the most significant phrases.

eg "New York" all docs that contain this phrase would be placed into a cluster. Finally these initial clusters are merged based on doc-doc similarity.

② Hierarchical Agglomerative clustering.

First the $N \times N$ doc similarity matrix is formed. Each doc is placed into its own cluster. The following two steps are repeated until one cluster exists.

The two clusters that have the highest similarity are found.

These two clusters are combined & the similarity between the newly formed cluster & the remaining clusters recomputed.

As the larger cluster is formed the clusters that merged together are tracked & form a hierarchy.

Assume docs A B C D E exist & a doc-doc sim matrix exists. At this pt each doc is a cluster by itself.

$\{ \{A\} \{B\} \{C\} \{D\} \{E\} \}$

$\{ \{A, B\} \{C\} \{D\} \{E\} \}$

After certain steps

$\{A, B, C, D, E\}$

However the history of the formation of this