

Improving IR performance from OCREd text using cooccurrence *RISOT 2012*

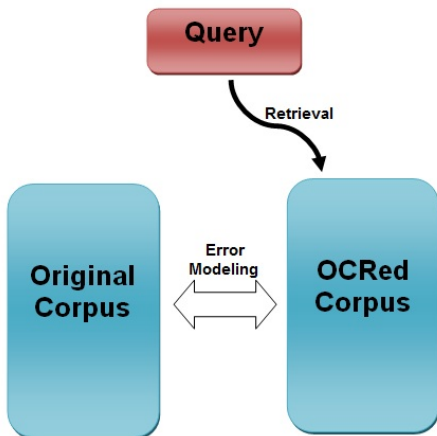
Kripabandhu Ghosh
and
Anirban Chakraborty
Indian Statistical Institute
Kolkata, India

Improving IR performance from OCRed text using cooccurrence RISOT 2012

Key Terms

- Co-occurrence : We say that two words co-occur if they appear in a window of certain number of words of each other in a document
- LCS_similarity : **LCS** stands for **L**ongest **C**ommon **S**ubsequence
 - $LCS(\textit{industry}, \textit{industrial}) = \textit{industr}$
 - $LCS_similarity(\textit{industry}, \textit{industrial}) = \frac{|LCS(\textit{industry}, \textit{industrial})|}{\max(|\textit{industry}|, |\textit{industrial}|)} = 0.7$

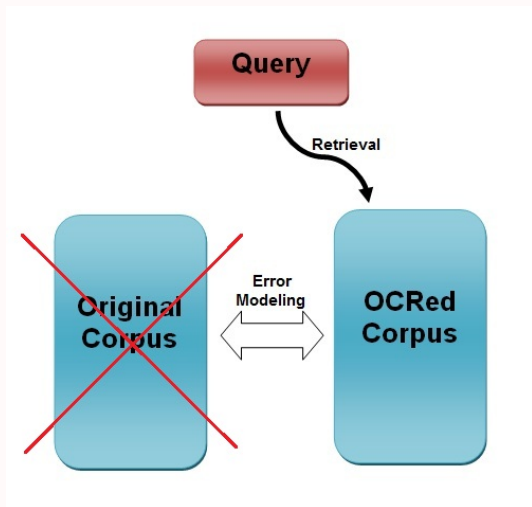
RISOT task



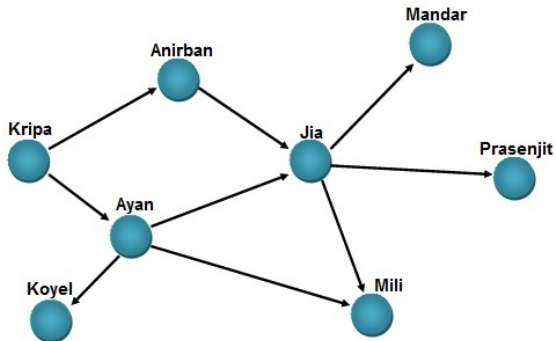
OCRed Resources

- Legal documents as hard copies
 - IIT CDIP 1.0 corpus - TREC Legal Ad Hoc
- SIGIR Digital Museum - Cleverdon, Salton, Sparck Jones

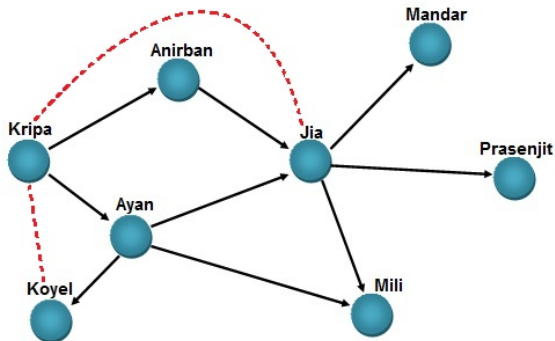
RISOT task - Without Original Corpus



Social Networks : Direct Connections



Social Networks : Indirect Connections



CLUSTERING

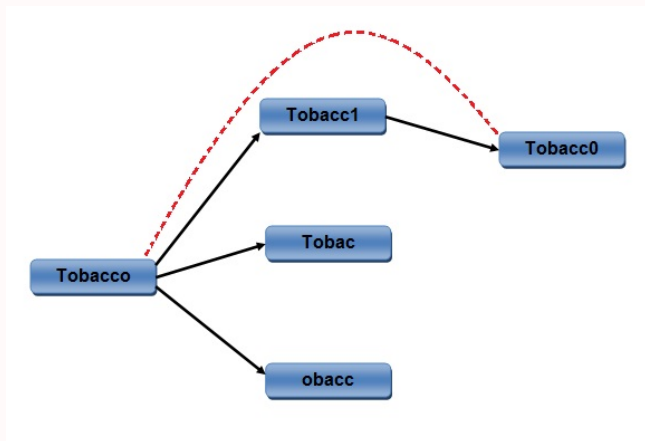
Clustering Algorithm (Phase I)

For each word w in the OCR'd corpus:

Let S_{w1} and S_{w2} be empty sets. Let $S_w = S_{w1} \cup S_{w2}$

- 1 For word w_1 co-occurring with w , calculate LCS_similarity between w and w_1 . Store w_1 in S_{w1} if $\text{LCS_similarity}(w, w_1) >$ some threshold T .
- 2 For each w' in S_{w1} , find the words w_2 co-occurring with w' such that $\text{LCS_similarity}(w, w_2) > T$. Include all these words in S_{w1} .
- 3 Repeat step (2) until no new word is added to S_{w1} .

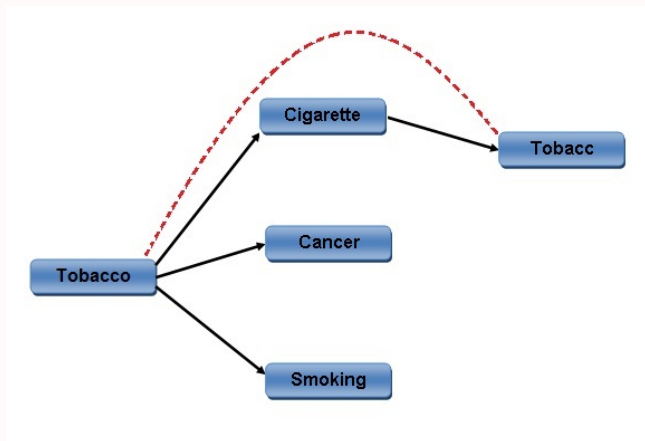
Clustering Algorithm (Phase I)



Clustering Algorithm (Phase II)

- 1 Consider top m (in terms of frequency in corpus) words co-occurring with w . For each such word w_3 , find the words w_4 cooccurring with w_3 such that $\text{LCS_similarity}(w, w_4) > T$. Include all these words in S_{w2} .
- 2 For each w'' in S_{w2} , find the words w_5 co-occurring with w'' such that $\text{LCS_similarity}(w, w_5) > T$. Include all these words in S_{w2} .
- 3 Repeat step (2) until no new word is added to S_{w2} .

Clustering Algorithm (Phase II)



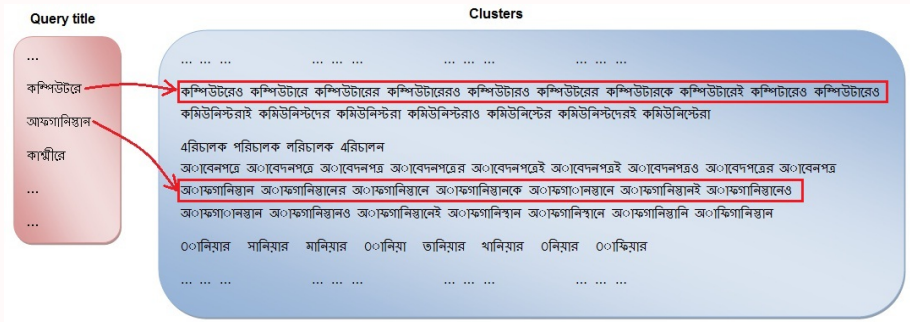
Query word cluster mapping

- 1 calculate $\text{LCS_similarity}(w_q, w_C)$, where w_C is a word in cluster C , for each word in the corpus
- 2 Choose all the clusters C for which $\text{LCS_similarity}(w_q, w_C)$ is greater than a high threshold
- 3 For each cluster C obtained from step (2) define $C' = C \cup \{w_q\}$
- 4 Create complete-linkage clusters from each C' of step (3) and keep those clusters containing w_q

Query word cluster mapping

- 5 For a cluster C , let us consider LCS_similarity between each pair of words in it. Let GM_C denote the Geometric Mean of LCS_similarity of all the pairs. Then, compute GM_C for each cluster given by step (4)
- 6 Select the cluster C with maximum GM_C as the appropriate cluster for w_q

Query word clusters

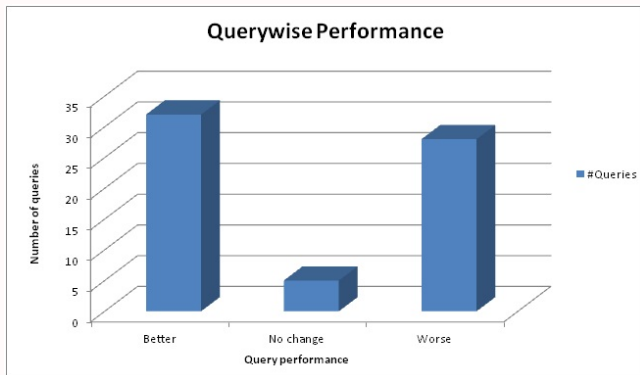


Results

Run	MAP	P5
Original text	0.2567	0.3485
OCRed text (baseline)	0.1791	0.2738
Proposed method on OCRed text	0.1974 ¹	0.2831

¹Not significant

Querywise Performance



Failure Analysis : clusters

The Good

... ..

কম্পিউটারেও কম্পিউটারে কম্পিউটারের কম্পিউটারেরও কম্পিউটারও কম্পিউটারের কম্পিউটারকে কম্পিউটারেই কম্পিটারেও কম্পিউটারেও
কমিউনিস্টরাই কমিউনিস্টদের কমিউনিস্টরা কমিউনিস্টরাও কমিউনিস্টের কমিউনিস্টদেরই কমিউনিস্টেরা
আবেদনপত্রে আবেদনপত্রে আবেদনপত্র আবেদনপত্রের আবেদনপত্রেই আবেদনপত্রই আবেদনপত্রও আবেদনপত্রের আবেদনপত্র
আফগানিস্তান আফগানিস্তানের আফগানিস্তানে আফগানিস্তানকে আফগানিস্তানে আফগানিস্তানই আফগানিস্তানেও আফগানিস্তান
আফগানিস্তানও আফগানিস্তানেই আফগানিস্তান আফগানিস্তানে আফগানিস্তানি আফগানিস্তান

... ..

and...

The Bad

... ..

4রিচালক পরিচালক লরিচালক 4রিচালন
0ানিয়ার সানিয়ার মানিয়ার 0ানিয়া তানিয়ার থানিয়ার 0নিয়ার 0াফিয়ার

... ..

Failure Analysis

- Compact clusters over all-inclusive clusters
- Re-clustering based on string match and co-occurrence
- Chance co-occurrence - harmful
- Incorporation of co-occurrence frequencies - essential

Brighter side

- Practical utility
- Language independent
- Context information - reliable
- Captures both erroneous and inflectional variants (effect of stemming)

THANK YOU