

# Stopword Removal: Why Bother? A Case Study on Verbose Queries

## ABSTRACT

Stopword removal has traditionally been an integral step in information retrieval preprocessing. In this paper we question the utility of this step for verbose queries on standard datasets. We show in seven FIRE test collections in four languages – Bangla, Hindi, Gujarati and Marathi, that stopwords removal does not lead to noticeable difference in retrieval performance as opposed to not removing stopwords. However, for European languages like English (TREC678 Ad Hoc) and French (CLEF 2005 to 2007), stopwords removal leads to a statistically significant drop in performance.

## CCS CONCEPTS

•Information systems →Information retrieval;

## KEYWORDS

Stopword removal, retrieval, test collections

### ACM Reference format:

. 2017. Stopword Removal: Why Bother? A Case Study on Verbose Queries. In *Proceedings of ACM CIKM conference, Pan Pacific, Singapore, November 2017 (CIKM'17)*, 2 pages.  
DOI:

---

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

*CIKM'17, Pan Pacific, Singapore*

© 2017 Copyright held by the owner/author(s). ...\$15.00

DOI:

Dataset	Number of documents	Number of topics
FIRE 2010 Ad Hoc Bangla	123047	50
FIRE 2011 Ad Hoc Bangla	377104	50
FIRE 2012 Ad Hoc Bangla	377111	50
FIRE 2010 Ad Hoc Hindi	149482	50
FIRE 2011 Ad Hoc Hindi	331599	50
FIRE 2011 Ad Hoc Gujarati	313163	50
FIRE 2010 Ad Hoc Marathi	99275	50
TREC678 Ad Hoc English	528155	150
CLEF 2005 to 2007 Ad Hoc French	177452	148

Table 1: Datasets.

Language	Source
Bangla	<a href="http://www.isical.ac.in/~fire/data/stopwords_list_ben.txt">http://www.isical.ac.in/~fire/data/stopwords_list_ben.txt</a>
Hindi	<a href="http://www.isical.ac.in/~fire/data/stopwords_list_hin.txt">http://www.isical.ac.in/~fire/data/stopwords_list_hin.txt</a>
Gujarati	<a href="http://irlab.daiict.ac.in/downloads/gujarati_stop_words.zip">http://irlab.daiict.ac.in/downloads/gujarati_stop_words.zip</a>
Marathi	<a href="http://members.unine.ch/jacques.savoy/clef/marathiST.txt">http://members.unine.ch/jacques.savoy/clef/marathiST.txt</a>
English	<a href="http://www.lemurproject.org/stopwords/stoplist.dft">http://www.lemurproject.org/stopwords/stoplist.dft</a>
French	<a href="http://members.unine.ch/jacques.savoy/clef/frenchST.txt">http://members.unine.ch/jacques.savoy/clef/frenchST.txt</a>

Table 2: Stopword sources.

Dataset	Language	With stopwords	Stopwords removed
FIRE 2010	Bangla	<b>0.4337</b>	0.4334
FIRE 2011	Bangla	<b>0.2910</b>	0.2898
FIRE 2012	Bangla	0.2465	<b>0.2472</b>
FIRE 2010	Hindi	<b>0.4695</b>	0.4569
FIRE 2011	Hindi	0.1639	<b>0.1642</b>
FIRE 2011	Gujarati	<b>0.2818</b>	0.2753
FIRE 2010	Marathi	<b>0.2587</b>	0.2552
TREC678	English	<b>0.2150*</b>	0.1971
CLEF 2005 to 2007	French	<b>0.2739*</b>	0.2581

Table 3: This table reports the retrieval performance in terms of Mean Average Precision. For each dataset the better performance value is shown in bold font. \* indicates that the performance difference is statistically significant ( $p < 0.05$ ) at 95% confidence interval.