Stop Word Lists in Document Retrieval Using Latent Semantic Indexing: an Evaluation

A N K Zaman School of Computer Science University of Guelph Guelph, ON, Canada azaman@uoguelph.ca



ABSTRACT: Removing stop words is very useful for many text processing applications e.g. text/document retrieval, cross language translation, text categorization, text summarization etc. In this world, different language has different stop word lists, and those are useful for text processing applications. Literature claims that the use of such lists improves retrieval performance. The goal of this research is to evaluate the effect of using English stop word lists in Latent Semantic Indexing (LSI)-based information retrieval (IR) systems with large text dataset. Here, three different lists are compared: two were compiled by IR groups at the University of Glasgow, and the University of Tennessee, and the third one is our own list developed at the University of Northern British Columbia. We also examined the case where stop words were not removed from the input dataset. Our research finds that using tailored stop word lists improves retrieval performance. On the other hand, using arbitrary (non-tailored) lists or not using any list reduces the retrieval performance of LSI-based IR systems with large text dataset.

Keywords: Document retrieval, Latent Semantic Indexing, Stop word list, Recall-precision

Received: 1 October 2011, Revised 6 November 2011, Accepted 12 November 2011

© 2012 DLINE. All rights reserved

1. Introduction

While working with text data, some particular words are very common in every document, and those have very little influence to distinguish documents from each other; these words are considered as stop words. In 1990, Deerwester et al., [1] removed a set of 439 common (stop) words from MED and CISI datasets. From our knowledge, this is the first stop word list used to work with latent semantic indexing (LSI). Example of stop words include articles (e.g., a, an, the), prepositions (e.g., at, by, in, to, from, with), conjunctions (e.g., and, but, as, because) etc.

2. Related Works

Proper stop word identification and removal are a central problem for many text processing applications in different domains. Stop words of a given text dataset might not be stop words of other datasets [2]. Stop words have significant impact on the text retrieval processes in different languages. In [3], Dolamic et al., evaluated two stop word lists with lengths 571 and 9 respectively on the Cross-Language Evaluation Forum (CLEF) dataset. According to the authors, the Divergence from Randomness (DFR) model shows lower retrieval performance when a short or no stop word list is removed from the input dataset. However, in case

of a revised Okapi (information retrieval system) implementation, retrieval performance does not show any significant difference whether a short, long or no stop word list is removed from the dataset. The authors also draw the same conclusion for other natural languages such as French, Hindi, and Persian. The main weaknesses of their work are that they use arbitrary stop word lists and do not use the *tf-idf* (term frequencyinverse document frequency) weighting scheme [4]. Zou et al., [5] show that the removal of stop words from Chinese text is important for Chinese word segmentation and improves the performance of Chinese text retrieval. The removal of stop words also improves systems' performance in Arabic IR and Arabic text summarization [6] [7]. The removal of stop words has positive impact on English text categorization [8]. Stop word removal also improves retrieval performance in case of cross-language IR. A number of crosslanguage based IR systems are reported in literature, e.g., Bengali-Hindi, Turkish-English, Japanese-English [9] [10] [11]. Schuemie et al., [12] removed stop words for cross-language IR for biomedical literature. The removal of stop words plays an important role in different text processing domains in different languages.

Although the Text REtrieval Conference (TREC) encourages text retrieval research, it does not provide any evidence or rule to use stop words in IR research. Different IR research groups use different stop word lists, and the size of these lists vary. As there is no standard stop word list for English text, one open question is the following: What are the effects of tailored (based on a certain dataset) vs. arbitrary (not tailored) stop word lists on LSI-based text retrieval systems with large datasets? In this study, we used TREC-8 LA Times dataset for our experiments. Two existing, arbitrary stop word lists, as well as our own tailored stop word list are considered (see table 2). The retrieval results are compared by removing three different sets of stop words from the input dataset (one at a time) and also without removing any stop words. The rest of the paper is organized as follows. Section 3 gives an overview of the LSI technique. Section 4 presents some characteristics of the input dataset and describes the way our stop word list is compiled. The experimental setup and retrieval results are presented in Section 5. Concluding remarks are in Section 6. Appendix A presents our developed stop word list (also called UNBC stop word list).

3. Latent Semantic Indexing

In 1988, Dumais et al. [13] introduced the idea of LSI to process textual data, and to handle synonymy, and polysemy. After that, LSI is used for many applications including document retrieval.

3.1 Overview

LSI is a method that exploits the idea of vector space model and singular value decomposition (SVD). SVD is an effective dimensional reduction scheme. It has been proved to be a very good choice for uncovering latent semantic structure [1]. SVD can be applied with an arbitrary rectangle matrix with the entries on the rows and columns. The matrix is then decomposed into three matrices containing singular vectors and/or singular values. These three matrices with special forms show a breakdown of the original matrix into linearly independent components or factors. Many of these components are very small, leading to an approximate model that contains many fewer dimensions. Thus, for IR purposes, SVD provides a reduced model for representing the *term-to-term*, *document-to-document* and *term-to-document* relationships. By dimension reduction, it is possible for documents with somewhat different profiles of term usage to be mapped into the same vector of factor values. This property helps to eliminate the noise in the original data, thus improving the reliability of the algorithm. Suppose we obtained a *td* term-by-document matrix *M* from the collection indexing process of the traditional vector space method. We can apply SVD on *M*, which is then decomposed into three special matrices *U*, *S* and *V*. The decomposition can be written as:

$$M = USV^{T}$$
(1)

U is the t * t orthogonal matrix (UU^T = I_t) having the left singular vectors of M as its columns, and V is the d * d orthogonal matrix (VV^T = I_d) having the right singular vectors as its columns, and S is the t * d diagonal matrix having the singular values $\sigma_1 \ge \sigma_2$ $\ge \cdots \ge \sigma_{\min(t,d)}$ of M in order along its diagonal. It should be noted that for any arbitrary matrix, such a factorization exists [14].

Generally, in (1), the matrices U, S and V must all be of full rank. However, SVD offers a simple strategy for optimal approximation to fit using smaller matrices [1]. If the singular values in S are ordered by size, the first k largest values may be kept and the remaining smaller ones set to zero. The product of the resulting matrices is a matrix M_k which is only approximately equal to M, and is of rank k. Since zeros were introduced into S, the representation can be simplified by deleting the zero rows and columns of S to obtain a new diagonal matrix S_k , and then deleting the corresponding columns of U and V to obtain U_k and V_k respectively. The rank-k model with the best possible least-squares-fit to M can be written as follows:

$$M_k = U_k S_k V_k^T \tag{2}$$

Where, M_k is a matrix of size t * d, U_k is of size t * k, S_k is of size k * k, and V_k is of size k * d.

SVD provides an optimal solution to dimensionality reduction in that it derives an orthonormal space, where the dimensions are ordered. Therefore, projecting the set of documents onto the k lowest dimensions is guaranteed to have, among all possible projections to a k dimensional space, the lowest possible least-square distance to the original documents.

3.2 Weighting Schemes

One significant issue in LSI-based IR systems is the term weighting, i.e., assigning weight to a term so that the assigned weight properly reflects the contribution of the term in distinguishing the considered document from other documents. Let, L_{ij} be the local weight of the term *i* in the document *j* and tf_{ij} be the frequency with which the term *i* appears in the document *j*. The local weight in terms of raw term frequency is defined as follows:

Raw term frequency:
$$L_{ii} = tf_{ii}$$
 (3a)

Let G_i be the global weight of the term *i*, let tf_i be the frequency of the term *i* in the entire collection, let df_i be the frequency of documents in which *i* occurs, and let *d* be the number of documents in the whole collection. The following equations define the *idf* and *tf-idf* weighting schemes:

$$idf: = G_i = \log\left(\frac{d}{df_i}\right) \tag{3b}$$

$$tf - idf : G_i = tf_{ij} \times idf_i \tag{3c}$$

For our experiments we used *tf-idf* weighting scheme to create term-document matrix.

4. Input Dataset and Stop Word Lists

The input data/text collection includes the articles published by the Los Angeles Times in the two year period from Jan 1, 1989 - December 31, 1990. Each file contains the articles from one day (e.g., a file with the name "LA123190" contains articles published on 31 Dec 1990). Every such file contains a number of documents (e.g., the LA123190 contains 134 different documents). Table 1 presents the important characteristics of the TREC-8 LA Times dataset.

Number of documents	131,321
Size of the input dataset	476MB
Average vocabulary size (approximately)	500
Average document size (approximately)	40 KB
Largest file size	828 KB (LA052089_0101)
Smallest size	352 Bytes (LA070189_000)
Number of words in the smallest file	91
Number of words in the largest file	167,045
Number of relevant files (out of 131,312 files) with respect to TREC-8	
query set	1,151

Table 1. Characteristics of the TREC-8 LA Times Dataset (1989, 1990)

Manning et al., [15] describe a way in their book to prepare a list of stop words: "the general strategy for determining a stop list is to sort the terms by collection frequency (the total number of times each term appears in the document collection), and then to take the most frequent terms, often hand-filtered for their semantic content relative to the domain of the documents being

Journal of E- Technology Volume 3 Number 1 February 2012

indexed, as a stop list, the members of which are then discarded during indexing." Our own stop word list has been compiled by following the above idea. It includes the University of Glasgow (319 words) and the University of Tennessee (439 words) stop word lists, 730 TREC file names (input dataset), 22 tag names (e.g., doc, docno, etc) and other words (e.g., alphanumeric words, roman numbers). Its total length is 1891 unique words. The algorithmic steps to create this stop word list are given below:

- Consider the terms whose frequency is at least 2 (a term must be present in the document at least twice).
- Create an initial stop word list by combining the stop word lists of the IR University of Tennessee and
- University of Glasgow groups (without duplication of terms in the list).
- Remove all the punctuation from the input TREC-8 LA Times dataset.

• Create a list of terms from the input dataset, in descending order of term frequencies, i.e., the term with the highest term frequency will be at the top of the list.

• Manually extract the special items to be added to the initial list (those terms are not already in the initial list) to create an extended stop word list.

- Add all file names to the initial list as every file contains file names, e.g., LA123190.
- Add all tag names, e.g., doc, docno, to the initial stop word list.
- Add roman numbers to the initial list, e.g., xvii.
- Add scale units, e.g., ft, mm, etc.
- Add adjectives and adverbs, e.g., ago.
- Add prefixes from words, e.g., non (as in non-governmental).
- Add special words, e.g., haven (as in haven't), doesn (as in doesn't).
- Add dates, e.g., 19th.
- Add foreign words as dataset in newspaper articles. Add suspicious words, e.g., aaftink, aachen, ora.
- Add other words, e.g., ext (telephone extension), 19th, z90, v6 (engine).

To compile this stop word list, first we searched the high frequency, low frequency, and then other special terms out of 132,785 terms in the frequency table. We repeat this in a number of cycles by removing different stop words from the TREC-8 LA Times dataset. Searching stop words is very time consuming as the dataset as well as the number of terms are large. The most difficult thing is to choose a word as a stop word. Since the TREC-8 LA Times dataset contains newspaper articles (on politics, sports, geography, history, science-technology, etc.), there are variations in the contents. Some characteristics of this dataset are presented in the table 1.

5. Experiments

This study finds out the effects of "*stop words/common words*" on a LSI-based text document retrieval process for the TREC-8 LA Times dataset. Evidence is developed to indicate the most effective stop word lists for LSI-based ad hoc IR processes for the mentioned dataset in table 1. We performed our experiment by removing the stop word lists mentioned in table 2, as well as without removing them. We applied Porter's stemming [16] to find the root words from the input text. 50 TREC-8 queries were used (associated with the mentioned dataset) to evaluate the retrieval performance. Our findings are presented in terms of recallprecision graph [14].

Table 3 shows the 10-point interpolated precision of the four different retrieval systems: the UNBC system which uses our stop word list developed at the University of Northern British Columbia; the UTen system which uses the University of Tennessee stop word list; the UGla system which uses the University of Glasgow stop word list; and the STEM system which does not use any stop word list. The recall-precision graph based on the results in Table 3 is shown in Figure 1.

5.1 Recall-Precision Graph

In Table 3, the recall value of 0.1 represents the top 10% of the retrieved documents (in the collection) which are relevant to a query set. As an example, using the UNBC system, the precision associated with the top 10% of the documents is 0.1757 (i.e.,

Dataset	Stop Word Lists	Stemming	Weighting Scheme	Number of Queries
	University of Glasgow			
TREC-8LA	University of Tennessee	Porter's	.6.116	7 0
Times	University of Northern British Columbia	Stemming	tf - 1df	50
	No list			

Table 2. Parameters for the study of stop words

10-Point Recall	UNBC	UTen	UGla	STEM
0.1	0.1757	0.1513	0.1221	0.1385
0.2	0.1108	0.0994	0.0864	0.1058
0.3	0.0888	0.0735	0.0799	0.0841
0.4	0.0724	0.0693	0.0722	0.0743
0.5	0.0678	0.0672	0.0694	0.0710
0.6	0.0659	0.0647	0.0671	0.0662
0.7	0.0646	0.0609	0.0660	0.0632
0.8	0.0621	0.0588	0.0628	0.0592
0.9	0.0554	0.0489	0.0578	0.0553
1.0	0.0436	0.0316	0.0374	0.0247

Table 3. 10 – point Interpolated Precision of the Four Systems



Figure 1. Recall-Precision Graph

17.57%). This value is calculated by interpolating the precision values of all 50 queries used for this research at the standard recall value 0.1.

The retrieval systems are compared in terms of precision in different standard recall points, e.g., 0.1, 0.2. For example, at recall point 0.3 (top ranked 30% documents), the precision values for UNBC is 8.88%, and it is 7.35% for UT. So, UNBC shows 1.53% (8.88%-7.35%) better retrieval performance than UT for the top 30% retrieved documents. If we look at the recall point 0.3 in figure 1, we can see the differences.

In the end, the system UNBC with extended stop word list provides the best result when compared to the three other systems. For the top 10% retrieval, it shows 5.37% better retrieval performance than UG, 3.68% better retrieval performance than STEM, and 2.44% better retrieval performance than UT. However, after top 40% retrieval all the systems show almost the same retrieval performance. Note that in STEM, we just applied Porter's stemming without removing stop words, and the retrieval performance is 1.64% better than UG's. From the above results, it is clear that the use of an arbitrary set of stop words reduces retrieval performance in case of LSI-based ah hoc IR with large dataset. Dai and Sun [17] showed in their paper that topic specific stop word removal from news stories helps to identify accurate event. This idea also proves our idea that tailored stop word list removal from specific text dataset improves document retrieval. Removal of proper stop word lists not only benefitted LSI based text processing but also text processing with other methods e.g. kK-Nearest neighbour(k-NN) [18], tagged Huffman compressed text searching [19] etc.

6. Conclusion

To identify and use of proper stop word lists from a specific dataset is very important in doing research with text. In our research, we investigated the performance of LSI by using three different stop word lists, and also, without using any stop word list, i.e., without removing stop words from the input dataset. Our main finding is that for a LSI-based ad hoc IR system, the use of an arbitrary stop word list reduces retrieval performance; for better retrieval performance, a tailored stop word list must be assembled for every unique large dataset.

Appendix A

a	a, a1, a10, a11, a12, a13, a14, a15, a16, a18, a19, a2, a20, a21, a22, a23, a24, a25, a26, a27, a28, a29, a3, a30, a300, a301, a31, a310, a32, a320, a320s, a321, a33, a330, a330s, a340, a340s, a35, a36, a4, a41, a43, a5, a6, a7, a8, a9, aa, aaa, aaaa, aaaah, aable, aachen, aad, aaf, aaftink, aagla, aah, aahed, aahing, aahs, aai, aar, aas, ab, aba, abb, abc, abl, about, above, abt, ac4, aca, acc, accordingly, ace, ach, acn, acoc, across, ad, ada, adl, adm, adn, ae, afc, afi, afl, after, afterwards, ag, again, against, agn, ago, ah, ahl, ahm, ak, al, ala, alf, all, allows, alm, almost, alone, along, alp, already, also, although, always, am, ama, amc, ami, amo, among, amongst, amp, an, ana, anc, and, ang, anh, ani, ann, another, any, anyawi, anybody, anyhow, anyone, anything, anyway, anywhere, ap, apart, appear, appropriate, ar, are, around, as, asa, ash, aside, aso, associated, ast, asu, at, att, av, ava, available, aw, away, awfully
b	b, b1, b10, b100, b12, b2, b210, b3, b4, b5, b52, b6, b7, b747, b8, b9, ba, ba2, ba3, ba4, ba6, baa, baa1, baa2, baa3, back, bbdo, be, became, because, become, becomes, becoming, been, before, beforehand, behind, being, below, ben, ber, beside, besides, best, better, between, beyond, blm, bmg, bnd, bo, boa, bol, boo, both, bp, brief, bs, bsd, bta, btu, bu, but, bv, by, bye, byline, byu
с	c, c1, c10, c11, c12, c13, c14, c15, c16, c17, c18, c19, c2, c20, c21, c22, c23, c24, c2h2, c3, c4, c5, c6, c7, c759915, c8, c9, cac, came, can, cannot, cant, cause, causes, cb, cbn, cbo, cc, ccaa, ccdc, cch, cct, ce, cee, cellrule, certain, cfa, ch, cha, changes, chi, chj, chp, chr, chu, cio, ck, clo, cmc, cmdr, cmv, co, column, come, consequently, contain, containing, contains, coq, correction, corresponding, cot, could, cpa, cpl, cra, crc, currently, cvj, cwl, cy, cya
d	d, d1, d10, d11, d12, d13, d14, d15, d16, d2, d3, d4, d5, d6, d7, d8, d9, date, dateline, day, db, de, dea, ded, def, described, dh, di, dib, did, didn, different, dl, dmc, do, doc, doc, doc, docid, docno, docno, docno, does, doesn, doing, don, done, doo, dosen, down, downwards, dr, ds, du, during, dy

e	e, e1, e5, each, eb, ebb, ec, ed, edit, ee, eel, eg, eh, ei, eight, eip, either, ek, el, elk, else, elsewhere, em, en, enough, eo, ep, epa, epp, er, es, et, etc, ev, even, ever, every, everybody, everyone, everything, everywhere, ew, ex, example, except, ext
f	f, f1, faa, far, fcc, fe, few, fg, fi, fifth, fig, first, five, fl, flo, followed, following, foo, for, former, formerly, forth, four, from, ft, fu, further, furthermore
g	g, g3, g4, g5, g6, ga, gc, ge, get, gets, given, gives, gn, go, gone, goo, good, got, gq, graphic, graphics, great, grp, gte, gto, gtp, gu
h	h, h1, h11, h13, h18, h2, h20, h2a, h2o, h3, h4, h5, h6, h7, ha, haa, had, hadn, hardly, has, have, haven, having, he, headline, hem, hence, her, here, hereafter, hereby, herein, hereupon, hers, herself, hi, him, himself, his, hither, hj, hmmm, hmmmmm, ho, hoo, how, howbeit, however, hoy, ht, hu, hy
i	i, ibf, iby, ic, ida, ie, if, ignored, ii, iii, ik, il, im, ima, immediate, in, ina, inasmuch, inc, indeed, indicate, ing, inner, insofar, instead, into, inward, ip, iq, ir, iri, irk, is, isl, isn, it, its, itself, iv, ivo, ix
j	j, j1, j10, j2, j3, j4, j5, j6, j8, ja, jc, ji, jo, jr, just
k	k, k1, k1n, k2, k2r, k9, ka, ka7, kan, kao, kee, keep, kept, kg2, kg7, kg8, kh1, kh2, kh7, kh8, know, knx, ko, ky
1	l, la, la2, la23, la89, la90, last, latter, latterly, le, least, length, less, lest, let, li, life, like, little, ll, lo, long, loo, lot, lou, lp, ls, lt, ltd, lx, ly
m	m, m1, m16, m2, m3, m4, m5, m6, m62, m71, m78, ma, maa, made, mae, make, mal, man, many, may, mc, mca, mcb, md, mdc, me, mea, meanwhile, mee, men, mfg, mg, mi, might, mk, mm, mme, mo, moc, moo, more, moreover, most, mostly, mott, mr, ms, msg, mt, mu, much, must, mvp, mx, my, myself
n	n, n1, n2, n5, n6, na, na3, na6, name, namely, ncols, ncr, ne, near, necessary, neither, nev, never, nevertheless, new, next, nfc, nfl, ng, ni, nine, nl, nlrb, nmb, no, nobody, nom, non, none, noone, nor, normally, not, nothing, nov, novel, now, nowhere, nr, nu
0	o, oat, oc, occ, och, od, oda, odd, of, off, oft, often, oh, oj, ol, old, on, once, one, ones, only, onto, ooh, ooo, ooooo, oooz, op, opt, or, ora, ord, ot, other, others, otherwise, ou, ought, our, ours, ourselves, out, outside, over, overall, ow, own
р	p, p1, p2, pa, page, part, particular, particularly, pb, pba, pc, pcb, pcc, pcl, pcp, pe, people, per, perhaps, pfc, pg, ph, pic, pj, placed, plc, please, plo, plus, possible, pp, pr, prc, pre, pro, probably, provides, pt, pta, ptl
q	q, q106, qb, qd2, qd3, qd4, qd7, qd8, qe2, qe4, qe5, qe6, qe7, qe8, qed, qf2, qf5, qf7, qg4, qg5, qg6, qh4, qi, qtr, que, quite
r	r, r1, r2, r2d2, r3, ra, ra3, ra6, rather, rb, rbi, rc, rda, re, really, relatively, respectively, right, rk, ro, roo, rowrule, rv
S	s, s1, s1630, s1w, s2, s358, s4, s605, said, same, sb, sba, sbk, sc, sca, scc, scca, sce, scr, sdg, sdy, se, second, secondly, section, see, seem, seemed, seeming, seems, self, selves, sensible, sent, serious, seven, several, shall, she, shoo, should, si, since, six, smc, smu, so, some, somebody, somehow, someone, something, sometime, sometimes, somewhat, somewhere, soo sou, specified, specify, specifying, spn, sq, sr, ss, ssi, ssy, st, state, still, stu, su, sub, subject, such, sup, sw, sy, syd
t	t, t4, ta, table, tablecel, tablecell, tablerow, take, taken, tcu, td, te, text, than, that, the, their, theirs, them, themselves, then, thence, there, thereafter, thereby, therefore, therein, thereupon, these, they, thi, third, this, thorough, thoroughly, those, though, three, through, throughout,

u u, u2, u60, ua, uc, ul, ulf, un, unc, und, under, unless, uno, until, unto, up, upa, upon, us, use, used, useful, uses, usf, using, usually, ut, uw v v, v20, v2500, v6, v8, v8v, v9l, v9w, va, value, various, ve, very, vh, vi, via, vii, viii, viil, vin, viz, vs, vt, uu, vw, vy w w, w6, w8, wa, was, wasn, way, wbc, wcb, wdm, we, wee, well, went, were, what, whatever, whatnot, when, whence, whenever, where, whereafter, whereas, whereby, wherein, whereupon, wherever, whether, which, while, whither, who, whoever, whole, whom, whose, why, wig, will, with, within, without, wk, word, words, work, world, would, wouldn, wow, wr, wt, wu, wuz, wwii x x, x2, x6, xi, xiii, xiv, xix, xr, xt xtra, xv, xvi, xxi, xxiii, xxiv, xxv, xxvi, xxx, xyz y y, y95, yap, ye, year, years, yet, yo, yoo, you, your, yours, yourself, yourselves, yve z z, z28, z90, zac, zero, zx, zz, zzzz, File la010189, la010190, la010289, la010290, la010389, la010390, la010489, la010490, la010589, la010990, la011089, la011090, la011789, la01190, la01189, la01190, la01189, la01190, la011289, la011290, la011389, la011390, la011489, la011490, la011589, la01190, la011289, la011290, la011389, la011390, la011489, la011490, la011289, la011290, la01289, la01290, la01289, la01290, la01289, la01290, la01289, la01290, la012789, la01290, la012789, la01290, la01289, la012990, la01289, la012990, la01289, la01290, la01389, la01390, la01389, la01390, la01489, la011490, la011589, la01299, la012990, la01289, la012990, la01289, la012990, la01289, la012990, la01289, la012990, la01289, la012990, la01289,		thru, thus, thy, time, tmc, tnt, to, together, too, tot, toward, towards, ts, tu, twa, twice, two, type
vv, v20, v2500, v6, v8, v8v, v9l, v9w, va, value, various, ve, very, vh, vi, via, vii, viii, viil, vin, viz, vs, vt, vu, vw, vyww, w6, w8, wa, was, wasn, way, wbc, wcb, wdm, we, wee, well, went, were, what, whatever, whatnot, when, whence, whenever, where, whereafter, whereas, whereby, wherein, whereupon, wherever, whether, which, while, whither, who, whoever, whole, whom, whose, why, wig, will, with, within, without, wk, wok, word, words, work, world, would, wouldn, wow, wr, wt, wu, wuz, wwiixx, x2, x6, xi, xiii, xiv, xix, xr, xt, xtra, xv, xvi, xxi, xxiii, xxiv, xxv, xxvii, xxx, xyzyy, y95, yap, ye, year, years, yet, yo, yoo, you, your, yours, yourself, yourselves, yvezz, z28, z90, zac, zero, zx, zz, zzzz,Filela010189, la010190, la010289, la010290, la010389, la010390, la010489, la010490, la010589, la010590, la010689, la010690, la010789, la010790, la011889, la011390, la011489, la011189, la011190, la011289, la011200, la011289, la011200, la011389, la011390, la011489, la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011889, la011390, la011489, la011490, la011589, la011200, la012589, la012190, la012289, la012290, la012389, la012390, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012990, la01389, la013000, la013189, la013190, la020189,	u	u, u2, u60, ua, uc, ul, ulf, un, unc, und, under, unless, uno, until, unto, up, upa, upon, us, use, used, useful, uses, usf, using, usually, ut, uw
ww, w6, w8, wa, was, wasn, way, wbc, wcb, wdm, we, wee, well, went, were, what, whatever, whatnot, when, whence, whenever, where, whereafter, whereas, whereby, wherein, whereupon, wherever, whether, which, while, whither, who, whoever, whole, whom, whose, why, wig, will, with, within, without, wk, wok, word, words, work, world, would, wouldn, wow, wr, wt, wu, wuz, wwiixx, x2, x6, xi, xiii, xiv, xix, xr, xt, xtra, xv, xvi, xxi, xxiii, xxiv, xxv, xxvii, xxx, xyzyy, y95, yap, ye, year, years, yet, yo, yoo, you, your, yours, yourself, yourselves, yvezz, z28, z90, zac, zero, zx, zz, zzzz,Filela010189, la010190, la010289, la010290, la010389, la010390, la010489, la010490, la010589, la010590, la010689, la010690, la010789, la010790, la010889, la010890, la010990, la011089, la011090, la011189, la011190, la011289, la011290, la011389, la011390, la011489, la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011389, la011390, la011489, la011490, la0112089, la012090, la012189, la011290, la012289, la012290, la012389, la011290, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012990, la013089, la013090, la013189, la013190, la020189, la013190, la01289, la012990, la013089, la013090, la013190, la02189, la013190, la021889, la012990, la013089, la013090, la013189, la013190, la02189, la01289, la012990, la012989, la012990, la013089, la013090, la013189, la013190, la02189, la013190, la01289, la012990, la013089, la013090, la013189, la013190, la020189,	v	v, v20, v2500, v6, v8, v8v, v9l, v9w, va, value, various, ve, very, vh, vi, via, vii, viii, vill, vin, viz, vs, vt, vu, vw, vy
x x, x2, x6, xi, xiii, xiv, xix, xr, xt, xtra, xv, xvi, xxi, xxiii, xxiv, xxv, xxvii, xxx, xyz y y, y95, yap, ye, year, years, yet, yo, yoo, you, your, yours, yourself, yourselves, yve z z, z28, z90, zac, zero, zx, zz, zzzz, File la010189, la010190, la010289, la010290, la010389, la010390, la010489, la010490, la010589, la010590, la010689, la010690, la010789, la010790, la010889, la010890, la010989, la010990, la011089, la011090, la011189, la011190, la011289, la011290, la011389, la011390, la011489, la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011389, la011390, la011489, la011490, la011589, la0112090, la012189, la011290, la012289, la012290, la012389, la012390, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012889, la012990, la012990, la013089, la013090, la013189, la013190, la020189,	W	w, w6, w8, wa, was, wasn, way, wbc, wcb, wdm, we, wee, well, went, were, what, whatever, whatnot, when, whence, whenever, where, whereafter, whereas, whereby, wherein, whereupon, wherever, whether, which, while, whither, who, whoever, whole, whom, whose, why, wig, will, with, within, without, wk, word, words, work, world, would, wouldn, wow, wr, wt, wu, wuz, wwii
yy, y95, yap, ye, year, years, yet, yo, yoo, you, your, yours, yourself, yourselves, yvezz, z28, z90, zac, zero, zx, zz, zzzz,Filela010189, la010190, la010289, la010290, la010389, la010390, la010489, la010490, la010589, la010590, la010689, la010690, la010789, la010790, la010889, la010890, la010989, la010990, la011089, la011090, la011189, la011190, la011289, la011290, la011389, la011390, la011489, la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011389, la011390, la011489, la011490, la011589, la011590, la012089, la012090, la012189, la012190, la012289, la012290, la012389, la012390, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012989, la012990, la013089, la013090, la013189, la013190, la020189,	X	x, x2, x6, xi, xiii, xiv, xix, xr, xt, xtra, xv, xvi, xxi, xxiii, xxiv, xxv, xxvii, xxx, xyz
z z, z28, z90, zac, zero, zx, zz, zzzz, File la010189, la010190, la010289, la010290, la010389, la010390, la010489, la010490, la010589, la010590, la010689, la010690, la010789, la010790, la010889, la010890, la010989, la010990, la011089, la011090, la011189, la011190, la011289, la011290, la011389, la011390, la011489, la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011389, la011390, la011489, la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011889, la011890, la011989, la011990, la012089, la012090, la012189, la012190, la012289, la012290, la012389, la012390, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012989, la012990, la013089, la013090, la013189, la013190, la020189,	у	y, y95, yap, ye, year, years, yet, yo, yoo, you, your, yours, yourself, yourselves, yve
Filela010189, la010190, la010289, la010290, la010389, la010390, la010489, la010490, la010589, la010590, la010689, la010690, la010789, la010790, la010889, la010890, la010989, la010990, la011089, la011090, la011189, la011190, la011289, la011290, la011389, la011390, la011489, la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011889, la011890, la011989, la011990, la012089, la012090, la012189, la012190, la012289, la012290, la012389, la012390, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012989, la012990, la013089, la013090, la013189, la013190, la020189,	Z	z, z28, z90, zac, zero, zx, zz, zzzz,
Names Ia010199, Ia010209, Ia010209, Ia010299, Ia010399, Ia010399, Ia010499, Ia010499, Ia010499, Ia010389, Names Ia010590, Ia010689, Ia010690, Ia010789, Ia010790, Ia010889, Ia010890, Ia010989, Ia010990, Ia011089, Ia011090, Ia011189, Ia011190, Ia011289, Ia011290, Ia011389, Ia011390, Ia011489, Ia011490, Ia011589, Ia011590, Ia011689, Ia011690, Ia011789, Ia011790, Ia011889, Ia011890, Ia011989, Ia011990, Ia012089, Ia012090, Ia012189, Ia012190, Ia012289, Ia012290, Ia012389, Ia012390, Ia012489, Ia012490, Ia012589, Ia012590, Ia012689, Ia012690, Ia012789, Ia012790, Ia012889, Ia012890, Ia012989, Ia012990, Ia013089, Ia013090, Ia013189, Ia013190, Ia020189,	File	
la011089, la011090, la011189, la011190, la011289, la011290, la011389, la011390, la011489, la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011889, la011890, la011989, la011990, la012089, la012090, la012189, la012190, la012289, la012290, la012389, la012390, la012489, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012989, la012990, la013089, la013090, la013189, la013190, la020189,	Names	la010139, la010190, la010239, la010239, la010399, la010399, la010439, la010439, la010439, la010439, la010309,
la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011889, la011890, la011989, la011990, la012089, la012090, la012189, la012190, la012289, la012290, la012389, la012390, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012989, la012990, la013089, la013090, la013189, la013190, la020189,	1 varies	la011089, la011090, la011189, la011190, la011289, la010290, la011389, la011390, la011489.
la011989, la011990, la012089, la012090, la012189, la012190, la012289, la012290, la012389, la012390, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012989, la012990, la013089, la013090, la013189, la013190, la020189,		la011490, la011589, la011590, la011689, la011690, la011789, la011790, la011889, la011890,
la012390, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790, la012889, la012890, la012989, la012990, la013089, la013090, la013189, la013190, la020189,		la011989, la011990, la012089, la012090, la012189, la012190, la012289, la012290, la012389,
la012889, la012890, la012989, la012990, la013089, la013090, la013189, la013190, la020189,		la012390, la012489, la012490, la012589, la012590, la012689, la012690, la012789, la012790,
		la012889, la012890, la012989, la012990, la013089, la013090, la013189, la013190, la020189,
1a020190, 1a020289, 1a020290, 1a020389, 1a020390, 1a020489, 1a020490, 1a020589, 1a020590,		la020190, la020289, la020290, la020389, la020390, la020489, la020490, la020589, la020590,
la020689, la020690, la020789, la020790, la020889, la020890, la020989, la020990, la021089,		1a020689, 1a020690, 1a020789, 1a020790, 1a020889, 1a020890, 1a020989, 1a020990, 1a021089,
la021090, la021189, la021190, la021289, la021290, la021389, la021390, la021489, la021490,		la021090, la021189, la021190, la021289, la021290, la021389, la021390, la021489, la021490,
la021589, la021590, la021689, la021690, la021789, la021790, la021889, la021890, la021989,		la021589, la021590, la021689, la021690, la021789, la021790, la021889, la021890, la021989,
la021990, la022089, la022090, la022189, la022190, la022289, la022290, la022389, la022390,		la021990, la022089, la022090, la022189, la022190, la022289, la022290, la022389, la022390,
la022489, la022490, la022589, la022590, la022689, la022690, la022789, la022790, la022889,		la022489, la022490, la022589, la022590, la022689, la022690, la022789, la022790, la022889,
la022890, la030189, la030190, la030289, la030290, la030389, la030390, la030489, la030490,		la022890, la030189, la030190, la030289, la030290, la030389, la030390, la030489, la030490,
la030589, la030590, la030689, la030690, la030789, la030790, la030889, la030890, la030989,		la030589, la030590, la030689, la030690, la030789, la030790, la030889, la030890, la030989,
la030990, la031089, la031090, la031189, la031190, la031289, la031290, la031389, la031390,		la030990, la031089, la031090, la031189, la031190, la031289, la031290, la031389, la031390,
la031489, la031490, la031589, la031590, la031689, la031690, la031789, la031790, la031889,		la031489, la031490, la031589, la031590, la031689, la031690, la031789, la031790, la031889,
la031890, la031989, la031990, la032089, la032090, la032189, la032190, la032289, la032290,		la031890, la031989, la031990, la032089, la032090, la032189, la032190, la032289, la032290,
la032389, la032390, la032489, la032490, la032589, la032590, la032689, la032690, la032789,		1a032389, 1a032390, 1a032489, 1a032490, 1a032589, 1a032590, 1a032689, 1a032690, 1a032789,
la032790, la032889, la032890, la032989, la032990, la033089, la033090, la033189, la033190,		la032790, la032889, la032890, la032989, la032990, la033089, la033090, la033189, la033190,
la040189, la040190, la040289, la040290, la040389, la040390, la040489, la040490, la040589,		la040189, la040190, la040289, la040290, la040389, la040390, la040489, la040490, la040589,
la040590, la040689, la040690, la040789, la040790, la040889, la040890, la040989, la040990,		la040590, la040689, la040690, la040789, la040790, la040889, la040890, la040989, la040990,
la041089, la041090, la041189, la041190, la041289, la041290, la041389, la041390, la041489,		la041089, la041090, la041189, la041190, la041289, la041290, la041389, la041390, la041489,
la041490, la041589, la041590, la041689, la041690, la041789, la041790, la041889, la041890,		la041490, la041589, la041590, la041689, la041690, la041789, la041790, la041889, la041890,
la041989, la041990, la042089, la042090, la042189, la042190, la042289, la042290, la042289,		la041989, la041990, la042089, la042090, la042189, la042190, la042289, la042290, la042289,
1a042390, 1a042489, 1a042490, 1a042589, 1a042590, 1a042689, 1a042690, 1a042789, 1a042790,		1a042390, 1a042489, 1a042490, 1a042589, 1a042590, 1a042689, 1a042690, 1a042789, 1a042790,
1042889, 10042890, 10042989, 10042990, 10043089, 10043090, 10050189, 10050190, 10050289,		1a042889, 1a042890, 1a042989, 1a042990, 1a043089, 1a043090, 1a050189, 1a050190, 1a050289,
12050290, 1205059, 12050590, 12050900, 12050900, 12050000, 12051090, 1205090, 1205090, 12051090, 1205090, 1205090, 1205090, 1205090, 1205090, 12051090, 12051090, 12051090, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 12050900, 1205090000000000000000000000000000000000		120507290, 12050505, 12050590, 12050469, 120504990, 120500509, 12050590, 12050590, 12050690, 12051020, 12051020
10000/09, 120000/90, 120000090, 12000090, 12000090, 12000090, 12000090, 1200000, 120000, 120000, 120000, 12000		10001/07, 120001/20, 12000007, 12000207, 12000207, 120002990, 120001089, 120001090, 120001090, 120001090, 12000
10031170,10031207,10031270,10031307,10031307,10031390,10031407,10031490,10031307,10031307,10031307,10031307,100		1a051170, 1a051207, 1a051270, 1a051507, 1a051570, 1a051407, 1a051490, 1a051507, 1a051590, 1a051590, 1a051500, 1a051500, 1a051500, 1a051500, 1a051500, 1a051500, 1a051500, 1a051500, 1a05150, 1a051500, 1a05
1a051007, 1a051070, 1a051707, 1a051770, 1a051007, 1a051070, 1a051970, 1a051970, 1a051970, 1a05200, 1a0520, 1a0500, 1a0500, 1a0500, 1a0500, 1a0500, 1a0500, 1a0500, 1a0500, 1a0500, 1		1a051007, 1a051070, 1a051707, 1a051770, 1a051007, 1a051070, 1a051970, 1a051970, 1a051970, 1a05200, 1a05200, 1a05200, 1a05200, 1a05200, 1a05200, 1a05200, 1a05200, 1a0520, 1a05200, 1a05200, 1a05200, 1a0520, 1a0520, 1a05200,
1a052090, 1a052109, 1a052190, 1a052209, 1a052290, 1a052290, 1a052590, 1a052590, 1a052499, 1a052490, 1a052800 1a052800 1a052800 1a052080		1a052579, 1a052197, 1a052197, 1a052297, 1a052297, 1a052577, 1a052579, 1a052497, 1a052497, 1a052497, 1a052497, 1a052497, 1a0522497, 1a052589, 1a05289, 1a052889, 1a05289, 1a052889, 1a05888, 1a058888, 1a058888, 1a058888, 1a058888, 1a05888, 1a058888, 1a058888, 1a0588888, 1a0588888, 1a0588888, 1a05888888888888888888888888888888888888
1a052507, 1a052570, 1a052007, 1a052070, 1a052787, 1a052779, 1a052679, 1a052677, 1a052777, 1a0527777, 1a057777, 1a057777, 1a057777, 1a057777, 1a057777, 1a057777, 1a0577777, 1a0577777, 1a0577777, 1a057777777, 1a05777777, 1a05777777, 1a05777777, 1a05777777, 1a0577777, 1a05777777777777, 1a0577777		1a052509, 1a052509, 1a052009, 1a052009, 1a052709, 1a052799, 1a052089, 1a052099, 1a05209, 1a052099, 1a05209, 1a052099, 1a05209, 1a0500000000000000000000000000000000000
1a060389 1a060390 1a060489 1a060490 1a060589 1a060590 1a060689 1a060689 1a060689		1a060389 1a060390 1a060489 1a060490 1a060589 1a060590 1a060689 1a060690 1a060789
1a060790, 1a060889, 1a060890, 1a060989, 1a060990, 1a061089, 1a061090, 1a061189, 1a061190.		la060790, la060889, la060890, la060989, la060990, la061089, la061090, la061189, la061190.

	la061289, la061290, la061389, la061390, la061489, la061490, la061589, la061590, la061689,
	1a061690, 1a061789, 1a061790, 1a061889, 1a061890, 1a061989, 1a061990, 1a062089, 1a062090,
	1a062189, 1a062190, 1a062289, 1a062290, 1a062389, 1a062390, 1a062489, 1a062490, 1a062589,
	1a062590 1a062689 1a062690 1a062789 1a062790 1a062889 1a062890 1a062989 1a062990
	1a063089 1a063090 1a070189 1a070190 1a070289 1a070290 1a070389 1a070390 1a070489
	1a070490 1a070589 1a070590 1a070689 1a070690 1a070789 1a070790 1a070889 1a070890
	$1_{a}070989$ $1_{a}070990$ $1_{a}071089$ $1_{a}071090$ $1_{a}071189$ $1_{a}071190$ $1_{a}071289$ $1_{a}071290$ $1_{a}071389$
	la071390 la071489 la071490 la071589 la071590 la071689 la071690 la071789 la071790
	la071390, la071490, la071990, la071390, la071390, la071090, la071090, la071790, la071790, la071790, la072289
	1a072200 1a072380 1a072390 1a072780 1a072790 1a072580 1a072590 1a072590 1a072680 1a072600
	1_{2}
	1_{2}
	1_{2}
	1_{2}
	$1_0001400, 1_0001000, 1_0001100, 1_0001100, 1_0001200, 1_0001200, 1_0001200, 1_0001200, 1_0001300, 1_000000, 1_00000, 1_000000, 1_00000, 1_00000, 1_00000, 1_000000, 1_000000, 1_000000, 1_000000, 1_000000, 1_000000, 1_0000000, 1_0000000, 1_00000000, 1_0000000000$
	1_{0}
	10001070, 10001707, 10001770, 10002007, 10002070, 10002107, 10002170, 10002207, 10002270,
	12002309, 12002390, 12002409, 12002490, 12002309, 12002390, 12002009, 12002090, 12002709, 120027000000000000000000000000000000000
	12002/90, 12002009, 12002090, 12002909, 12002990, 12003009, 12003090, 12003109, 12003190,
	10000500 10000200 10000200 10000200 10000700 10000200 10000200 10000200 10000000
	1a090390, 1a090089, 1a090090, 1a090789, 1a090790, 1a09089, 1a090890, 1a090890, 1a090989, 1a090990,
	120014001, 12001500, 12001100, 12001100, 12001200, 12001200, 12001200, 12001200, 12001200, 12001200, 12001200,
	12001080 1-001000 1-002080 1-002000 1-002180 1-002180 1-002280 1-002280 1-002280
	12091989, 12091990, 12092089, 12092090, 12092189, 12092190, 12092289, 12092290, 12092389, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 12092898, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289, 1209289
	12092390, 12092489, 12092490, 12092389, 12092390, 12092689, 12092690, 12092789, 12092790,
	12092889, 12092890, 12092989, 12092990, 12093089, 12093090, 12100189, 12100190, 12100289,
	la100290, la100389, la100390, la100489, la100490, la100589, la100590, la100689, la100690,
	la100/89, la100/90, la100889, la100890, la100989, la100990, la101089, la101090, la101189,
	la101190, la101289, la101290, la101389, la101390, la101489, la101490, la101589, la101590,
	la101689, la101690, la101789, la101790, la101889, la101890, la101989, la101990, la102089,
	la102090, la102189, la102190, la102289, la102290, la102389, la102390, la102489, la102490,
	la102589, la102590, la102689, la102690, la102789, la102790, la102889, la102890, la102989,
	la102990, la103089, la103090, la103189, la103190, la110189, la110190, la110289, la110290,
	la110389, la110390, la110489, la110490, la110589, la110590, la110689, la110690, la110789,
	la110790, la110889, la110890, la110989, la110990, la111089, la111090, la111189, la111190,
	la111289, la111290, la111389, la111390, la111489, la111490, la111589, la111590, la111689,
	la111690, la111789, la111790, la111889, la111890, la111989, la111990, la112089, la112090,
	la112189, la112190, la112289, la112290, la112389, la112390, la112489, la112490, la112589,
	la112590, la112689, la112690, la112789, la112790, la112889, la112890, la112989, la112990,
	la113089, la113090, la120189, la120190, la120289, la120290, la120389, la120390, la120489,
	la120490, la120589, la120590, la120689, la120690, la120789, la120790, la120889, la120890,
	la120989, la120990, la121089, la121090, la121189, la121190, la121289, la121290, la121389,
	la121390, la121489, la121490, la121589, la121590, la121689, la121690, la121789, la121790,
	la121889, la121890, la121989, la121990, la122089, la122090, la122189, la122190, la122289,
	la122290, la122389, la122390, la122489, la122490, la122589, la122590, la122689, la122690,
	la122789, la122790, la122889, la122890, la122989, la122990, la123089, la123090, la123189,
	la123190
Dates	19th
Duito	1/11

References

[1] Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., Harshman, R. A. (1990). Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41 (6) 391-407.

Journal of E- Technology Volume 3 Number 1 February 2012

[2] Dragut, E. C., Fang, F., Sistla, A. P., Yu, C. T., Meng, W. (2009). Stop word and related problems in web interface integration. The Proceedings of the Very Large Database Endowment, PVLDB 2 (1), p. 349–360.

[3] Dolamic, L., Savoy, J. (2010). When stopword lists make the difference. Journal of the American Society for Information Science and Technology, 61: p. 200–203.

[4] Zaman, A. N. K., Brown, C. G. (2010). Latent semantic indexing and large dataset: Study of termweighting schemes. Fifth International Conference on Digital Information Management (ICDIM), p.1-4.

[5] Zou, F., Wang, F. L., Deng, X., Han, S. (2006). Evaluation of Stop Word List in Chinese Language. International Conference on Language Resources and Evaluation (LREC), Genoa, Italy.

[6] El-Khiar., I. (2006). Effects of Stops Words Elimination for Arabic Information Retrieval: A Comparative Study. *In*: International *Journal of Computing and Information Sciences*, 4 (3).

[7] Azmi, A., Al-thanyyan, S. (2009). Ikhtasir - A user selected compression ratio Arabic text summarization system. International Conference on Natural Language Processing and Knowledge Engineering, p.1-7.

[8] Feng Xia, Tian Jicun, Liu Zhihui. (2009). A Text Categorization Method Based on Local Document Frequency. Sixth International Conference on Fuzzy Systems and Knowledge Discovery, p.468-471.

[9] Mandal, D., Gupta, M., Dandapat, S., Banerjee, P., Sarkar, S. (2008). Bengali and Hindi to English CLIR evaluation. Advances in Multilingual and Multimodal Information Retrieval (8th Workshop of the Cross-Language Evaluation Forum). Budapest, Hungary Springer Verlag, p. 95–102.

[10] Celebi, E., Sen, B., Gunel, B. (2009). Turkish - English cross language information retrieval using LSI. 24th International Symposium on Computer and Information Sciences, p.634-638.

[11] Li and Shawe-Taylor, Li, Y., Shawe-Taylor, Y. (2006). Using KCCA for Japanese–English crosslanguage information retrieval and document classification. *Journal of Intelligent Information System* 27 (2) 117–133.

[12] Dumais, S. T., Furnas, G. W., Landauer, T. K., Deerwester, S. (1988). Using latent semantic analysis to improve information retrieval. Conference on Human Factors in Computing, New York: ACM, p.281-285.

[13] Schuemie M., Trieschnigg D., Kraaij, W. (2007). Cross Language Information Retrieval for Biomedical Literature. Sixteenth Text REtrieval Conference, Gaithersburg, MD, USA.

[14] Schutze, H., Silverstein, C. (1997). Projections for efficient document clustering. In: ACM/SIGIR(SIGIR Forum 31), p. 74-81.

[15] Christopher D. Manning, Prabhakar Raghavan, Hinrich Schütze. (2008). An Introduction to Information Retrieval. Cambridge University Press Cambridge.

[16] Porter, M. (2002). An algorithm for suffix stripping, avaliable :http://tartarus.org/~martin/PorterStemmer/ Access time.

[17] Dai, Xiangying, Sun, Yunlian, (2010). Event identification within news topics. International Conference on Intelligent Computing and Integrated Systems, p.498-502.

[18] de Oliveira Gomes, N., Passos, E. P. L. (2011). Text categorization study case: Patents' application documents. Sixth IEEE Conference on Industrial Electronics and Applications, p.446-450.

[19] Edleno Silva de Moura, Gonzalo Navarro, Nivio Ziviani, Ricardo Baeza-Yates. (2000). Fast and flexible word searching on compressed text, ACM Trans. Inf. Syst. 18, p.113-139.