DEFENCE **R&D** DÉFENSE

# Latent Semantic Analysis (LSA) tools

Natalia Derbentseva
Peter J. Kwantes
Philip Terhaar

Canadä

# Latent Semantic Analysis (LSA) tools

Natalia Derbentseva
Peter J. Kwantes
Philip Terhaar

## Defence R&D Canada – Toronto

Principal Author

*Original signed by Natalia Derbentseva*

...............................................................................................................

Natalia Derbentseva

Defence Scientist

Approved by

*Original signed by Keith Stewart*

...............................................................................................................

Keith Stewart

Head of Socio-Cognitive Systems Section

# Table of contents

This page intentionally left blank.

# 1    Introduction

Latent Semantic Analysis (LSA; Landauer & Dumais, 1997; Landauer, Foltz, & Laham, 1998) is a computational model that uses a large collection of unstructured documents to construct semantic representations for words. The representations are based on a statistical analysis of the terms' occurrences within and across documents, and take the form of a vector. The semantic similarity between resulting word representations can be compared by calculating their cosine. After an LSA space is created, it can be queried to provide a word-to-word, word-to-document, or document-to-document comparisons to determine their semantic similarity by taking the cosine of the angle formed by their vector representations. When comparing words to documents (and sometimes documents to each other), LSA uses what is called a "bag of words" approach to representing the semantic contents of a document. In a "bag of words", the order of terms in a document does not matter, and the semantic representation of a document is formed by summing the vectors of all its content words.

LSA requires a relatively large set of short documents to generate a semantic space (i.e., from several hundreds to tens of thousands). The large number helps to ensure that many words from a variety of contexts (operationally defined as documents) are available for analysis.

Before a semantic space for the terms in the corpus is constructed by our version of LSA, the documents need to be prepared. Punctuation is removed from the documents, words are brought to lower case, and words that appear in most or all documents in the corpus, such as articles, prepositions, and others that do not help to differentiate the contextual uses of the terms are removed from the documents, so called "stop words". Typically, the same set of stop words is removed from every document in a corpus before analysis. However, under certain conditions different stop word lists have to be removed from different groups of documents before they can be combined into a single corpus and analyzed. Because LSA usually requires several thousands of documents, removing custom stop word lists from them is a time-consuming and labour intensive task. One of the tools reported in this Technical Note (TN) supports this task.

One of LSA's requirements is that for a word to be included in the space, it must occur in at least two documents. This poses a potential challenge when a corpus for LSA analysis needs to be created from a limited number of relatively short documents. In this situation, a large number of words can potentially be excluded from the analysis because they do not meet the minimum document occurrence criteria. At least partially, it occurs not because there are many unique words, but because the same word is used in different forms in different documents (for example, singular and plural forms of nouns). Thus, bringing different forms of the same word to a single form across documents in a training corpus can partially diminish their unnecessary exclusion from the corpus. The tools described in this TN implemented stemming of words prior to the construction of the corpus and the LSA space.

LSA builds its semantic representation from a consideration of words' contextual use. That is, words that tend to be used in similar contexts will tend to have a semantic association. Not surprisingly then, the same two words can have different representations depending on the collection of documents used to create a semantic space, For example, the word "mouse" means something different in the context of computers than in the context of animals or vermin. Hence, depending on the purpose of the analysis, the same set of words or texts can be analyzed with

different LSA spaces. The WIKIPEDIA subcorpora tool (Stone and Dennis, 2012) allows the user to create custom corpora from a Wikipedia archive.

Finally, the existing functionality of DRDC's SEMMOD package allows pair-wise comparisons of words or texts, one pair at a time with a single LSA space. Therefore, a large number of such comparisons are a time-consuming and a labour intensive task. One of the tools described in this TN provides automation of multiple LSA comparisons and supports application of several LSA spaces.

## 1.1    Utilities to Support Latent Semantic Analysis (LSA)

This TN documents a set of tools that were developed to support the ongoing work at DRDC Toronto that applies. The described set of tools builds on the WIKIPEDIA Subcorpora Tool (Stone & Dennis, 2012) and the Semantic Models (SEMMOD) module (v 1.5, Dennis & Stone, 2011) that were developed by the Ohio State University (OSU) under the contract W7711-067985/001/TOR for a Technology Investment Fund (TIF) project (15da05) and an Applied Research Project (ARP, 15ah) . The tools described in this TN extend the capabilities offered by the OSU modules and support 2 main functions:

1. Preparation of a custom corpus ready for LSA's semantic space generation. In addition to standard procedures that are performed on a collection of documents to prepare them for LSA space generation, such as removing punctuation and stop words, the tools developed also allow the following operations:

    ◆ The application of a customized stop word list to a single document file prior to its inclusion into a corpus;

    ◆ The stemming the words in a file using the Porter stemming algorithm;

    ◆ The exclusion of words that occur in fewer than a specified number of documents.

2. Automating the process of conducting document-to-document comparisons where documents are created from different LSA spaces, with the following options:

    ◆ application of a customized stop word list to the input file;

    ◆ stemming the words in the input file;

    ◆ specifying multiple LSA spaces to be applied to the input file;

    ◆ generating output in the form of a document-by-document table or a single column to facilitate subsequent analysis of the results.

●

# 2    LSA supporting tools

This section describes the tools that support LSA including their functions, input requirements, output formats, and execution requirements. The tools were developed in the Python 2.6.5 programming language for the Linux environment. In the interest of time no graphical user interface (GUI) was developed, and individual files containing scripts are run from the command line. The source code is currently stored on DRDC Toronto public server Pluto at the following location: \\Pluto\public\VISTIF\LSATools. Print outs of the source code are provided in Annexes A-C.

## 2.1    Tools to support custom corpus generation

This sub-section describes the tools that were developed to support custom corpus generation, including i) possibility to apply a custom stop words list to a single document prior to its inclusion into a corpus, ii) stemming the words in a file using the Porter stemming model; iii) excluding from the corpus words that occur in fewer than the specified number of documents.

### 2.1.1    Removing custom stop words from a file

This module was developed to allow the user to remove a custom set of words from a single text file before the file is included into a corpus file. This module is useful when different stop word lists have to be applied to different documents before they are combined into a single corpus. NOTE: prior to removing stop words, punctuation is also removed from the file.

*File:* *CustomStopW.py* (source code is in Appendix A)

*Command line prompt example:*

$    python    CustomStopW.py    -i    /home/MyDocuments/TextFiles/Agreeableness.txt    -s /home/MyDocuments/StopListFiles/Agreeableness_Stop.txt

*Options and arguments:*

> *-i*  followed by the input file name. Requires either a filename with a complete path, or just a name of a file in the default location: "*input/*" subfolder. If this option is omitted, the default input file will be processed, which is "*input/default.txt*".

> *-s*  followed by the stop words list file.  Requires either a filename with a complete path, or just a name of the file in the default location: "*CleaningFiles/*" subfolder. If this option is omitted, the default stop words list file will be used: "*CleaningFiles/stopList_Words.txt*".

*Input requirements:*

> ◆ Input file: Plain text, preferably in UTF-8 format with no special requirements. Default input file location is "*input/*" subfolder.

- Stop words list file: Plain text, preferably in UTF-8 format. Each stop word must appear on a separate line.

*Output:*

Punctuation and the words found in the stop words list file are removed from the input file, and the result is stored in a new file, which is saved in the same location as the input file. The line breaks in the input file are preserved. The name of the output file complies with the following naming convention:

<input file>_CSWr.<input file extension>.

*Execution requirements:*

The folder containing the *CustomStopW.py* must also contain the following files:
- *corpusCleaningTools.py*
- *EntityClassify.py*

The folder containing the stop word list file must also contain the punctuation stop list: "*stopList_Punctuation.txt*".

## 2.1.2    Custom corpus preparation

This module prepares a text document for LSA space generation, and outputs a collection of documents in a file with the extension *.cor*. It allows the use to apply a custom stop words list to the input file, stem the words, and specify the minimum frequency of word's occurrence.

*File:* *CustomCorpusPreparation.py* (source code is in Appendix B)

*Command line prompt example:*

$ python CustomCorpusPreparation.py -i /home/MyDocuments/TextFiles/Agreeableness.txt –f 3 -s /home/MyDocuments/StopListFiles/Agreeableness_Stop.txt –t

*Options and arguments:*

- *-i*  followed by the input file name. Requires either a filename with a complete path, or just a name of a file in the default location: "*input/*" subfolder. If this option is omitted, the default input file will be processed, which is "*input/default.txt*".

- *-f*  to indicate the minimum number of documents in which a word must occur to be included in the corpus. This argument must be followed by an argument. The default value is 2.

- *-s*  followed by the stop words list file.  Requires either a filename with a complete path, or just a name of the file in the default location: "*CleaningFiles/*" subfolder. If this

option is omitted, the default stop word list file will be used: "*CleaningFiles/stopList_Words.txt*".

-*t* no argument is required. If "-*t*" is included, the words will be stemmed after the documents are cleaned, and they will not be stemmed if "-*t*" is omitted. The default is to omit stemming.

***Input requirements:***

- ◆ Input file: Plain text, preferably in UTF-8 format with no special requirements. Default input file location is "*input/*" subfolder.

- ◆ Stop words list file: Plain text, preferably in UTF-8 format. Each stop word must appear on a separate line.

- ◆ The folder containing the stop word list file must also contain the punctuation stop list: "*stopList_Punctuation.txt*".

***Output:***

A corpus ready for LSA space generation: Punctuation and stop words removed, only words that appear with sufficient frequency across documents are retained. The documents are separated by a blank line. The result is stored in a new file which is placed in the same location as the input file with the following naming convention:

<input file>.cor

The words that were removed from the documents are stored in a file with the name

<input file>.cor.removed

***Execution requirements:***

- ◆ *nltk.stem.porter*

- ◆ *semmod.lsa*

- ◆ The folder containing the *CustomStopW.py* must also contain the following files:

  - ◆ *corpusCleaningTools.py*

  - ◆ *EntityClassify.py*

- ◆ The folder containing the stop word list file must also contain the punctuation stop list: "*stopList_Punctuation.txt*".

## 2.1.2.1   Word stemming module

Words can appear in several forms, like singular and plural versions of the nouns *book* and *books*. Different forms of a word are treated as different terms by LSA. As a result, one could argue that the number of unique words recognised by the system in a collection of documents is somewhat

inflated. Such inflation can pose a problem for generating a semantic space from a relatively small collection of short documents. In such a collection, a substantial number of words could be excluded from the analysis because of their "uniqueness" in the corpus, thus, resulting in a limited semantic space.

Word stemming can be used to mitigate this issue. Word stemming is a process that brings affixed forms of a word to its base form, that is, its stem. Reducing words to their stems decreases the number of unique words in a corpus, and increasing their frequency in the corpus. Such processing could improve the quality of a semantic space constructed from a relatively small number of short documents.

A function was developed that transforms a string into a collection of candidate stems. To reduce noise introduced by the stemming process itself, a validation step checks whether a stemmed form is a recognized word itself. If the stem is a recognized word, then the original word in the document is changed into its stemmed form. If the stemmed form is not recognized, the letters that were removed by stemming are added back to the stemmed form one by one. The check is repeated after each letter is added until a recognized word is found.

This function uses the Porter stemmer module from the Natural Language Toolkit (NLTK) package and it also uses an existing semantic space built using LSA to check for the words' existence. This function can be called from other modules by importing it and passing the required arguments to it.

 *Resides in file:* *CustomCorpusPreparation.py* (code is in Appendix B)

*Usage example:*

    from CustomCorpusPreparation import StemmerWithLSAcheck

    …

        StemmedDocuments = StemmerWithLSAcheck (ListOfDocuments, FileNameBase)

*Required arguments:*

- ◆ **ListOfDocuments** – either a list or a dictionary that contains strings to be stemmed
- ◆ **FileNameBase** – input file name without extension. It is used to create stemming output files

*Returns:* Stemmed strings in the form they were provided to the function, i.e., either a list or a dictionary

Also creates:

- ◆ a file with the record of stemming steps, file name  <inputfilename>_STemDebug" and
- ◆ a file that contains a list of the original words and their form after stemming, file name: <inputfilename>_AfterStemmingWords.txt"

*Execution requirements:*

All modules required by the *CustomCorpusPreparation.py* file (see 2.1.2), and the LSA semantic space contained in the file, *LSAspaces/tasaCleaned_fromPluto.lsa*

## 2.2 Tools supporting document-by-document comparisons using LSA with multiple semantic spaces

This module automates document-by-document comparisons, and it allows the user to measure the semantic similarity among all document pairs in an input file. Further, the comparisons can be conducted on representations derived from multiple semantic spaces from LSA at the same time. The module generates a document-by-document similarity matrix and saves it as comma separated values (CSV) in a separate ASCII file for each semantic space. It supports the following options:

- ◆ applying a customized stop words list to the input file;

- ◆ stemming the words in the input file;

- ◆ analysing the input file with multiple LSA spaces;

- ◆ generating output in the form of a document-by-document table , as well as a single column created by concatenating the columns of the matrix to facilitate subsequent analysis.

*File: doc_by_doc_Multi_LSA_Custom_STOPlist.py*  (source code is in Appendix C)

*Command line prompt example:*

$ python doc_by_doc_Multi_LSA_Custom_STOPlist.py -i /home/MyDocuments/TextFiles/Agreeableness.txt -f 3 -s /home/MyDocuments/StopListFiles/Agreeableness_Stop.txt -l /home/MyDocuments/LSAspaces -t

*Options and arguments:*

- *-i*  followed by the input file name. Requires either a filename with a complete path, or just a name of a file in the default location: "*input/*" subfolder. If this option is omitted, the default input file will be processed, which is "*input/default.txt*". The input file has to be properly formatted, see section "Input requirements" below for instructions.

- *-s*  followed by the stop words list file.  Requires either a filename with a complete path, or just a name of the file in the default location: "*CleaningFiles/*" subfolder. If this option is omitted, the default stop words list file will be used: "*CleaningFiles/stopList_Words.txt*".

*-l* followed by the name of the directory that contains all of the LSA spaces to be applied to the input file. If not specified, the default folder is "*LSAspaces/*".

*-t* no argument is required. If "-*t*" is included, the words will be stemmed after the docs are cleaned. They will not be stemmed if "-*t*" is omitted. The default is to omit stemming.

### Input requirements:

- ◆ Input file: All documents to be analyzed must be compiled into a single plain text file, in which each document is on a separate line; and each line (document) must begin with the document code, which will be used as the document identifier in all output files. Preferred format for the input file is UTF-8. Default input file location is "*input/*" subfolder.

- ◆ Stop words list file: Plain text, preferably in UTF-8 format. Each stop word must appear on a separate line.

- ◆ The folder containing the stop word list file must also contain the punctuation stop list: "*stopList_Punctuation.txt*".

### Output:

Two CSV files for each semantic space are created in the "*csv/*" sub-folder. One of the *.csv* files contains a document-by-document similarity matrix populated with cosine values above the upper diagonal from the given semantic space; the second *.csv* file (with the *_COLUMN* suffix) formats the same information into a single column. Files are named with the following convention:

<input-file name>_<StopList file name>_LSA_<name of the LSA space>.csv

<input-file name>_<StopList file name>_LSA_<name of the LSA space>_COLUMN.csv

### Execution requirements:

- ◆ *nltk.stem.porter*

- ◆ *semmod.lsa*

- ◆ *numpy*

- ◆ The folder containing the *doc_by_doc_Multi_LSA_Custom_STOPlist.py* must also contain the following files:

    - ◆ *corpusCleaningTools.py*

    - ◆ *EntityClassify.py*

    - ◆ *CustomCorpusPreparation.py (*contains the stemming function*)*

- ◆ The folder containing the stop word list file must also contain the punctuation stop list: "*stopList_Punctuation.txt*".

# 3    Concluding remarks

The tools described in this TN automate certain aspects of the otherwise labour-intensive and time-consuming process of pre-processing text for LSA and conducting multiple analyses. Although these tools were developed to address specific data analysis needs, the functionality that they support (e.g., document-by-document comparison) is fairly general. For example, the document-by-document comparison could be used in validating the semantic analysis component of the Analysis of Semantic and Social Networks (ASSN) tool.

The tools described in this TN will likely be developed further, given the nature of the ongoing work at DRDC Toronto. The purpose of this TN is to document the functionality developed up to date, to disseminate the availability of such functionality among DRDC Toronto colleagues who might benefit from them, and to reduce duplication of efforts in the future.

We expect to improve the flexibility and usability of these tools in the future, and to develop other functionality to support LSA and semantic analysis in general. DRDC Toronto can take the lead in developing a more comprehensive Python-based LSA toolkit, and make it available as an open source library to the general community of users interested in application of LSA. Such sharing could facilitate further co-development of the toolkit by the community.

This page intentionally left blank.

# Annex A  *CustomStopW.py* **Source Code**

```python
1   #!/usr/bin/env python
2
3   ################################################################################
4   #
5   #  Allows to apply a custom stop list to a file. Might be particularly useful when different
6   #  special sets of words must be removed from different files before they are combined in a
7   #  single corpus. Punctuation is also removed from the file.
8   #
9   #  ARGUMENTS:
10  #     -i for input file. Requires an argument. If path is not specified with
11  #        the file name, then the program will look for it in the "input/" sub-folder
12  #        relative to this script file.
13  #
14  #     -s for a stop words list file. If path is not indicated then the program will look for
15  #        it in the "CleaningFiles/" sub-folder. If stop list file is not specified or doesn't
16  #        exist, the program will use the "CleaningFiles/stopList_Words.txt".
17  #        NOTE: the foder with the stop words list file must also contain the
18  #                  "stopList_Punctuation.txt" file.
19  #
20  #  OUTPUT:
21  #     - File with the input file name and "-CSWr" suffix is saved in the input file's location.
22  #
23  ################################################################################
24
25
26  from corpusCleaningTools import removeStopWordsCaseInsensitive as RE_Stop
27  from corpusCleaningTools import replacePunctuationWithSpace as RE_Punct
28  import os, sys, getopt, pdb, glob
29
30  def main(argv):
31      INPUT_PATH = os.path.join(APP_PATH, 'input/')
32      FILES_PATH = os.path.join(APP_PATH, 'CleaningFiles/')
33      infilename = "default.txt"
34      stoplist = "stopList_Words.txt"
35
36      try:
37          opts, args = getopt.getopt(argv, "i:s:", ["input=", "stoplist="])
38      except getopt.GetoptError:
39          print "\nARGUMENT ERROR"
40          sys.exit(2)
41
42      for opt,arg in opts:
43          if opt in ("-i", "--input"):
44              #check if the input filename contains a path
45              if os.path.dirname(arg):
46                  INPUT_PATH = os.path.dirname(arg)
47                  infilename = os.path.basename(arg)
48              else:
49                  infilename = arg
50          elif opt in ("-s", "--stoplist"):
51              #check if the stoplist filename contains a path
52              if os.path.dirname(arg):
53                  FILE_PATH = os.path.dirname(arg)
54                  stoplist = os.path.basename(arg)
55              else:
56                  stoplist=arg
57
58      InputFile = os.path.join(INPUT_PATH, infilename)
59      StopListFile = os.path.join(FILE_PATH, stoplist)
60      puncFile = os.path.join(FILES_PATH, 'stopList_Punctuation.txt')
61
62      infile=open(InputFile,'r')
63      myText=infile.read()
64
65      print "\n\nApplying STOP LIST:     %s\nto the INPUT FILE:     %s\n" % (StopListFile,InputFile)
66      C_text = RE_Punct (myText, puncFile)
67      C_text = RE_Stop(C_text, StopListFile)
```

```python
68          NameBase, extent = os.path.splitext(InputFile)
69          outfilename="%s_CSWr%s" % (NameBase,extent)
70          CleanFile=open(outfilename, 'w')
71          CleanFile.write(C_text)
72
73          infile.close()
74          CleanFile.close()
75
76
77   if __name__ == "__main__":
78       main(sys.argv[1:])
```

# Annex B  *CustomCorpusPreparation.py* **Sourse Code**

```python
1   #!/usr/bin/env python
2
3   ###############################################################################
4   #
5   #     Takes an input file where each doc is on a single line.
6   #     User can specify the min number of docs in which each word has to
7   #     occur to be included in the corpus; and whether to stem the words or not.
8   #
9   #   ARGUMENTS:
10  #        -i for input file. Requires an argument. If path is not specified with
11  #           the file name, current location of this .py file is the default location.
12  #
13  #        -f to indicate the minimum number of docs in which a word must occur to
14  #           be included in the corpus. Must be followed by an argument.
15  #           The default value is 2.
16  #
17  #        -s for a stop words list file. If path is not indicated then the program will look for
18  #           it in the "CleaningFiles/" sub-folder. If stop list file is not specified or doesn't
19  #           exist, the program will use the "CleaningFiles/stopList_Words.txt".
20  #           NOTE: the foder with the stop words list file must also contain the
21  #                 "stopList_Punctuation.txt" file.
22  #
23  #
24  #        -t no argument is required. If "-t" is included, the words will be
25  #           stemmed after the docs are cleaned, and they will not be stemmed
26  #           if"-t" is omitted. The default is to omit stemming.
27  #
28  #     PROCESS:
29  #        - The input file is parsed and each line is treated as a separate document;
30  #        - Each document is cleaned with corpusCleaningTools.py;
31  #        - Words are stemmed (if -t is included in the command line);
32  #        - Words that occur in less than the specified number of docs are removed.
33  #     OUTPUT:
34  #        - File with .cor extension is saved in the input file's location.
35  #
36  #     Additional files generated:
37  #        - <inputfilename>.cor.removed lists the removed words;
38  #        - <inputfilename>_STemDebug" file records stemming steps;
39  #        - <inputfilename>_AfterStemmingWords.txt" contains the list of
40  #          original words and form after stemming.
41  #
42  ###############################################################################
43
44
45  import corpusCleaningTools as cleaner
46  import os, sys, getopt, pdb, glob
47
48  #from nltk.stem.lancaster import LancasterStemmer
49  from nltk.stem.porter import PorterStemmer
50  from semmod.lsa import lsa
51
52
53  def main(argv):
54      INPUT_PATH = os.path.join(sys.path[0], 'input/')
55      FILES_PATH = os.path.join(sys.path[0], 'CleaningFiles/')
56      infilename = 'default.txt'
57      stoplist= 'stopList_Words.txt'
58
59      XX=2
60      excludedDOCs=0
61      S=0
62
63      #Parsing options
64      try:
65          opts, args = getopt.getopt(argv, "i:f:s:t", ["inputfile=", "frequency=", "stoplist=",
    "stemming"])
66      except getopt.GetoptError:
```

```python
67              print "\nARGUMENT ERROR"
68              sys.exit(2)
69
70          for opt,arg in opts:
71              if opt in ("-i", "--input"):
72                  #check if the passed on input filename contains a path
73                  if os.path.dirname(arg):
74                      INPUT_PATH = os.path.dirname(arg)
75                      infilename = os.path.basename(arg)
76                  else:
77                      infilename = arg
78              elif opt in ("-f", "--frequency"):
79                  XX=int(arg)
80              elif opt in ("-s", "--stoplist"):
81                  if os.path.dirname(arg):                    #check if the specified stoplist filename
    contains a path
82                      FILES_PATH = os.path.dirname(arg)
83                      stoplist = os.path.basename(arg)
84                  else:
85                      stoplist = arg
86              elif opt in ("-t", "--stemming"):
87                  S=1
88          #End of option parsing
89
90          InputFile = os.path.join(INPUT_PATH, infilename)
91          stopFile = os.path.join(FILES_PATH, stoplist)
92          puncFile = os.path.join(FILES_PATH, 'stopList_Punctuation.txt')
93          NameBase, extent = os.path.splitext(InputFile)
94
95          #specify cleaning parameters and functions
96          myCleaningParams=[None,None,puncFile,None,None,None,stopFile]
97          MyCleanFunctions =
    ['lowerText','removeFormatting','replacePunctuationWithSpace','removeMultipleWhiteSpace','removeNumbers'
98
99          print "\n\n\nCUSTOM CORPUS PREPARATION RUNNING.........\n"
100         print "\n     Input file: %s\n" % InputFile
101         print "     Words that appear in less than %d documents will be removed\n" % XX
102
103         DOCS_C = [] #list for cleaned doc
104
105         infile = open(InputFile,'r')
106
107         for doc in infile: #go through lines in the input file, treat each line as a doc
108             if doc.strip(): #check if line is not empty
109                 #clean the doc with the corpusCleaningTools.py
110                 C_text=cleaner.runCleaningFunctions(doc, MyCleanFunctions, myCleaningParams)
111                 #Add the clreaned doc to the list of cleaned docs
112                 DOCS_C.append(C_text)
113
114         infile.close()
115
116         if S ==1:
117             print "   STEMMING with LSA CHECK ........."
118             DOCS_C = StemmerWithLSAcheck(DOCS_C, NameBase)
119         else:
120             print "\n   STEMMING WILL NOT BE PERFORMED!\n"
121
122         print "  Removing words that occur in less than %d documents:\n" % XX
123         WordsDocFrequency = WordsInDocsFrequency(DOCS_C)
124
125         DOCS_C =  RemoveLowFWords(WordsDocFrequency, DOCS_C, XX, NameBase)
126         print "\n   CREATING OUTPUT FILE:      %s.cor\n" % NameBase
127         #create an output .cor file
128         outfile=open("%s.cor" % NameBase,'w')
129         for doc in DOCS_C:
130             outfile.write ("%s\n\n" % doc)
131         print "CUSTOM CORPUS PREPARATION COMPLETE!\n\n"
```

```
132
133        outfile.close()
134
135
136
137   #************************************************************************************
138   #
139   #                              WORD FREQUENCY
140   #
141   #************************************************************************************
142        """
143        counts the number of documents in which each unique word appeared;
144        returns a dictionary of words with their doc frequency
145        """
146
147   def WordsInDocsFrequency(DOCS_C):
148
149        #create a dictionary that will have the number of documents in which each word from the
      corpus occured
150        words_in_docs_frequency = {}
151        words_doc_temp = [] #temp list of words to keep track of words in the same doc
152
153        for text in DOCS_C:
154            words_doc_temp[:]=[] #empty the temp list to get ready for a new doc
155            for word in text.split():
156                if word:
157                    if not word in words_doc_temp:
158                        words_doc_temp.append(word)
159                        if word in words_in_docs_frequency:
160                            words_in_docs_frequency[word]+=1
161                        else:
162                            words_in_docs_frequency[word]=1
163
164        return words_in_docs_frequency
165
166
167
168
169   #************************************************************************************
170   #
171   #     Remove words that occur in less than the specified
172   #     number of documents
173   #
174   #************************************************************************************
175
176   def RemoveLowFWords(WDFrequency, Docs, XX, NameBase):
177        #remove words that appear in less than XX number of documents
178
179        outfileremoved=open("%s.cor.removed" % NameBase,'w')
180
181        print "     Removed words are saved in: %s.cor.removed" % NameBase
182
183        removed_words={}
184        DOC_C_U=[]
185        i=0
186        excludedDOCs=0
187        for text in Docs:
188            string = ""
189            i+=1
190            for item in text.split(" "):
191                if item:
192                    if WDFrequency[item]>=XX:
193                        if string:
194                            string +=" %s" % item
195                        else:
196                            string=item
197                    else:
```

```python
198                         if item in removed_words:
199                             removed_words[item]+=1
200                         else:
201                             removed_words[item]=1
202
203                 if string:
204                     DOC_C_U.append(string)
205                 else:
206                     excludedDOCs+=1 #count excluded docs
207             for word, fr in removed_words.iteritems():
208                 outfileremoved.write("%s %d\n" % (word, fr))
209
210         print "\n        ...%d   UNIQUE words were found" % len(WDFrequency)
211         print "        ...%d   words were excluded because they occured in less than %d documents" % (len
        (removed_words), XX)
212         print "        ...%d   WORDS remain in the corpus" % (len(WDFrequency)-len(removed_words))
213         print "\n        ...%d   DOCUMENTS found in the INPUT file" % len(Docs)
214         print "\n        ...%d   DOCUMENTS are included in the final CORPUS" % len(DOC_C_U)
215         print "        ...%d   DOCUMENTS were EXCLUDED\n\n" % (excludedDOCs)
216
217         outfileremoved.close()
218
219         return DOC_C_U
220
221
222
223
224
225     #*********************************************************************************************
226     #
227     #                                   STEMMER with LSA
        check
228     #
229     #*********************************************************************************************
230
231     def StemmerWithLSAcheck(Docs, NameBase):
232         LSAname = os.path.join(sys.path[0], 'LSAspaces/tasaCleaned_fromPluto.lsa')
233         #LSAname = '/home/natalia/Devel/lsaplayground/Corpora/tasaCleaned_fromPluto.lsa'
234         #LSAname = '/home/natalia/Devel/lsaplayground/Corpora/WikiSubCorpRandom50000.lsa'
235
236         print "\n      The following LSA space will be used to check stemmed words' existance:
        \n           %s\n" % LSAname
237
238         if type(Docs)==type(dict()):
239             DOCS_C_S={}
240             for key,line in Docs.iteritems():
241                 DOCS_C_S[key]=stemAline(line, LSAname, NameBase)
242         elif type(Docs)==type(list()):
243             DOCS_C_S=[]
244             for line in Docs:
245                 DOCS_C_S.append(stemAline(line, LSAname, NameBase))
246         else:
247             print "\nError in StemmerWithLSAcheck function. A list or dictionary is required.\n"
248             sys.exit()
249         print "  STEMMING COMPLETE!\n"
250         return DOCS_C_S
251
252
253     #*********************************************************************************************
254     #
255     #     Stemms a line. Requires as inputs:
256     #           - line of text
257     #           - LSA space name
258     #           - NameBase (path plus the input filename)
259     #
260     #*********************************************************************************************
261
```

```python
def stemAline(line, LSAname, NameBase):
    Pst = PorterStemmer()
    stemfile=open("%s_STemDebug" % NameBase, 'a')
    stemmedwordsfile=open("%s_AfterStemmingWords.txt" % NameBase, 'a')
    stemfile.write("\nLINE:   %s\n" % line)
    LSAspace=lsa(LSAname) #lsa space to be used to check words' existance
    StemmedDoc=""

    for Original_Word in line.split():

        Stemmed_Word_Porter= Pst.stem(Original_Word)
        stemfile.write("\nOriginal WORD: __%s__         Stemmed WORD:  __%s__ " % (Original_Word,
Stemmed_Word_Porter))
        FoundWord=0
        try:
            IsVector=LSAspace.getTermVector(str(Stemmed_Word_Porter))
            #if vector exists then LSA knows the stemmed word and add the word to the doc string
            FoundWord=1
            stemfile.write("\nStemmed word found, adding %s\n" % Stemmed_Word_Porter)
            stemmedwordsfile.write("%s  -> %s\n" % (Original_Word, Stemmed_Word_Porter))
            if StemmedDoc:
                StemmedDoc +=" %s" % Stemmed_Word_Porter
            else:
                StemmedDoc=Stemmed_Word_Porter
        except KeyError:
            lenDif = len(Original_Word)-len(Stemmed_Word_Porter)
            if lenDif>0:
                for NextChar in range(1,lenDif+1):
                    if FoundWord ==0:
                        stemfile.write("\n  Trying ...  %s \n" % Original_Word[:(len
(Stemmed_Word_Porter)+NextChar)])
                        try:
                            IsVector=LSAspace.getTermVector(str(Original_Word[:(len(Stemmed_Word_Porter)
+NextChar)]))
                            #if vector exists then the stemmed+char(s) word exists and will add the
latest word form to the doc
                            FoundWord=1
                            stemfile.write("\nStemmed word + %d char found, adding %s\n" % (NextChar,
Original_Word[:(len(Stemmed_Word_Porter)+NextChar)]))
                            stemmedwordsfile.write("%s  ->  %s\n" % (Original_Word, Original_Word[:(len
(Stemmed_Word_Porter)+NextChar)]))
                            if StemmedDoc:
                                StemmedDoc+=" %s" % Original_Word[:(len(Stemmed_Word_Porter)+NextChar)]
                            else:
                                StemmedDoc= Original_Word[:(len(Stemmed_Word_Porter)+NextChar)]
                        except KeyError:
                            #the word not found, allow to add another char from the original
                            stemfile.write("\n   LSA DOESN'T KNOW  %s \n" % Original_Word[:(len
(Stemmed_Word_Porter)+NextChar)])
                            pass


            if FoundWord==0:
                stemfile.write("\nStemmed word NOT found, adding ORIGINAL %s\n" % Original_Word)
                stemmedwordsfile.write("%s  ->  %s\n" % (Original_Word, Original_Word))
                if StemmedDoc:
                    StemmedDoc+=" %s" % Original_Word
                else:
                    StemmedDoc=Original_Word
    stemfile.write("Stemmed line:   %s\n\n" % StemmedDoc)

    stemfile.close()
    stemmedwordsfile.close()

    return StemmedDoc

```

```
322   ###############################################################################
323
324
325
326   if __name__ == "__main__":
327       main(sys.argv[1:])
```

# Annex C   *doc_by_doc_Multi_LSA_Custom_STOPlist.py Sourse Code*

```python
1   #!/usr/bin/env python
2
3   ################################################################################
4   #
5   #    Document by Document LSA comparison:
6   #
7   #       Takes a text file, treats each line in the input file as a document and runs document by
8   #       document LSA comparison using all of the LSA spaces located in a specified folder.
9   #       Outputs one .csv file for each LSA space.
10  #
11  #    ARGUMENTS:
12  #      -i for input file. If path is not indicated then the program will look for it in the
13  #         "input/" sub-folder relative to this script file. If input file is not specified or
14  #         doesn't exist, the program will use the "input/default.txt".
15  #
16  #         INPUT FILE: A single text file that contains all the documets. Each document must be on a
17  #         single line and must begin with a document id followed by a period ("."). Line breaks are
18  #         treated as separators between documents. The program splits the inputfile into individual
19  #         documetns/lines using subject id as a key.
20  #
21  #      -s for a stop words list file. If path is not indicated then the program will look for
22  #         it in the "CleaningFiles/" sub-folder.
23  #         If stop list file is not specified or doesn't exist, the program will use the
24  #                          "CleaningFiles/stopList_Words.txt".
25  #
26  #         NOTE: the foder with the stop words list file must also contain the
27  #                  "stopList_Punctuation.txt" file.
28  #
29  #      -l for folder that contains LSA spaces. If not specified, the default folder is "LSAspaces/".
30  #         The program will identify all .lsa files in the folder and will run a line by line LSA
31  #         analysis using each of these spaces, creating an output file for each lsa space.
32  #
33  #      -t if included, the words will be stemmed after they are cleaned, and will not if
34  #         "-t" is omitted. Default is not to stem the words.
35  #
36  #
37  #  OUTPUT: Two CSV files for each LSA space are created in the "csv/" sub-folder. One of the .csv
38  #          files contains a doc by doc matrix populated with cosine values above the upper
        diagonal
39  #          from the given LSA space, while the second file (with the _COLUMN suffix) formats the
        same
40  #          information into a single column. Files are named with the following name
        convention:
41  #                          <input-file name>_<StopList file name>_<name of the LSA space>.csv
42  #                          <input-file name>_<StopList file name>_<name of the LSA space>_COLUMN.csv
43  #
44  #
45  #  Debug file: is saved in the "debug/Debug_<input file name>".
46  #
47  ################################################################################
48
49
50  print "\nRUNNING...\n"
51
52  #imports
53  #from pseudolizer import pseudolize
54  import corpusCleaningTools as cleaner
55  from semmod.lsa import lsa
56  from numpy import *
57  from numpy.linalg import *
58  from numpy.random import *
59  import os, sys, getopt, pdb, glob
60
61  from CustomCorpusPreparation import StemmerWithLSAcheck as stemm
62
63
64  def main(argv):
```

```python
65
66      APP_PATH = sys.path[0]
67      DEBUG_PATH = os.path.join(APP_PATH, 'debug/')
68      INPUT_PATH = os.path.join(APP_PATH, 'input/')
69      FILES_PATH = os.path.join(APP_PATH, 'CleaningFiles/')
70      OUTPUT_PATH = os.path.join(APP_PATH, 'csv/')
71      LSA_PATH = os.path.join(APP_PATH, 'LSAspaces/')
72      infilename = "default.txt"
73      stoplist = 'stopList_Words.txt'
74      S=0
75
76      try:
77          opts, args = getopt.getopt(argv, "hi:s:l:t", ["help", "inputfile=", "stoplist=","lsafolder=",
        "stemming"])
78      except getopt.GetoptError:
79          print "\nARGUMENT ERROR"
80          sys.exit(2)
81
82      for opt,arg in opts:
83          if opt in ("-h", "--help"):
84              print "\nDOC by DOC multi LSA analysis \n\nUse\n    -i to specify the input file,\n        the
        default location for the input file if path to the file is not specified is        /input/ sub-
        folder \n"
85              print "Use\n    -s to specify a stop words list file, the default location for cleaning files
        is:\n        /CleaningFiles/ sub-foledr\n"
86              print "Use\n    -l to specify a directory that contains LSA spaces (.lsa files).\n        The
        default directory is /LSAspaces/ sub-folder\n"
87              print "Use\n    -t to indicate that the words have to be stemed.\n        The words WILL NOT
        BE STEMMED if this option is omited\n\n"
88              sys.exit()
89          elif opt in ("-i", "--input"):
90              if os.path.dirname(arg):                        #check if the specified input filename contains
        a path
91                  INPUT_PATH = os.path.dirname(arg)
92                  infilename = os.path.basename(arg)
93              else:
94                  infilename = arg
95          elif opt in ("-s", "--stoplist"):
96              if os.path.dirname(arg):                        #check if the specified stoplist filename
        contains a path
97                  FILES_PATH = os.path.dirname(arg)
98                  stoplist = os.path.basename(arg)
99              else:
100                 stoplist = arg
101
102         elif opt in ("-l", "--lsafolder"):
103             if os.path.dirname(arg):                        #check if the specified stoplist filename
        contains a path
104                 if os.path.isabs(os.path.dirname(arg)):
105                     LSA_PATH = os.path.dirname(arg)
106         elif opt in ("-t", "--stemming"):
107             S=1
108
109
110
111     # define
112     InputFile = os.path.join(INPUT_PATH, infilename)
113     stopFile = os.path.join(FILES_PATH, stoplist)
114     puncFile = os.path.join(FILES_PATH, 'stopList_Punctuation.txt')
115     DebFileName = "%s_DEBUG" % infilename
116     Debugfile = os.path.join(DEBUG_PATH, DebFileName) #initiate debug file
117     NameBase, extent = os.path.splitext(InputFile)
118     LSAfiles = []                                      # list that will contain names of all .lsa
        files found in the LSA directory
119
120     if not os.path.isfile(InputFile):
121         print "\n\n???\n   EXITING: Could not find INPUT FILE %s\n   Specify a valid INPUT FILE with
```

```python
                 option -i and try again\n???\n\n" % InputFile
122          sys.exit()
123        if not os.path.isfile(stopFile):
124          print "\n\n???\n   EXITING: Could not find STOP LIST FILE %s\n   Specify a valid STOP LIST FILE
           with option -s and try again\n???\n\n" % stopFile
125          sys.exit()
126        if not os.path.isfile(puncFile):
127          print "\n\n???\n   EXITING: Could not find STOP LIST FILE %s\n   Make sure the file %s is in the
           %s directory, and try again\n???\n\n" % (puncFile, os.path.basename(puncFile), os.path.dirname
           (FILES_PATH))
128          sys.exit()
129        if not os.path.isdir(LSA_PATH):
130          print "\n\n???\n   EXITING: Could not find LSA directory %s\n   Specify a valid LSA directory
           with option -l and try again\n???\n\n" % LSA_PATH
131          sys.exit()
132
133
134        print "\n\nPREPARING FOR DOCUMENT BY DOCUMENT LSA ANALYSIS......\n"
135        print "\n...USING INPUT FILE: %s" % InputFile
136        print "...USING STOP LIST FILE: %s" % stopFile
137        if S==0:
138          print "...WORDS WILL NOT BE STEMMED"
139        elif S==1:
140          print "...WORDS WILL BE STEMMED"
141        print "\n...Folder\n %s\n   contains the following *.lsa files:\n" % LSA_PATH
142        for LSAfile in glob.glob(os.path.join(LSA_PATH, '*.lsa')):
143          LSAfiles.append(os.path.basename(LSAfile))
144          print "  %s" % os.path.basename(LSAfile)
145
146
147
148        #         Split the input file into individual lines (docs)
149
150        # Open the input file for reading
151        infile = open(InputFile,'r')
152        DebFile = open(Debugfile, 'w') #open debug file
153
154        myDOCS = {} # empty dictionary for docs
155
156        DebFile.write ("\n************************************************************\n")
157        DebFile.write ("**                        INPUT                        **\n")
158        DebFile.write ("************************************************************\n\n")
159        DebFile.write ("\nTHE INPUT FILE line by line: \n\n" )
160
161        print "\n PROCESSING the input file.....\n"
162        # Print the file line by line
163        print "   ...Splitting the input file...."
164        for line in infile:
165          DebFile.write (line) #write each line into the debug file
166          SubjectID, Para = line.split(".", 1) #split each line into Subject# and the doc
167          myDOCS[SubjectID]= Para.strip()  # create a dictionary of docs, with subject # as a key and
           strip front and back spaces...
168
169        # close input file
170        infile.close()
171
172        print "      Number of documents extracted:   %d\n" % len (myDOCS)
173        DebFile.write ("\n\nTOTAL NUMBER OF FILES IS: %d" % len (myDOCS))
174        DebFile.write ("\n\n\nTHE DICTIONARY OF DOCS CONTAINS THE FOLLOWING: \n" )
175
176        #specify cleaning parameters and functions
177        myCleaningParams=[None,None,puncFile,None,None,None,stopFile]
178        MyCleanFunctions =
           ['lowerText','removeFormatting','replacePunctuationWithSpace','removeMultipleWhiteSpace','removeNumbers'
179
180        # Clean each doc/line with corpusCleaningTools.py
181        print "   ...Cleaning the documents with corpusCleaningTools.py with the following options:\n
```

```python
      * lowerText\n       * removeFormatting\n    * replacePunctuationWithSpace\n       *
      removeMultipleWhiteSpace\n    * removeNumbers\n    * removeSingleCharacters\n     *
      removeStopWordsCaseInsensitive\n"
182    for n, t in myDOCS.iteritems():
183        DebFile.write ("\nSubject # %s:\nORIGINAL TEXT: %s" % (n,t))
184        if t.strip(): #check if line is not empty
185            myDOCS[n]= cleaner.runCleaningFunctions(t, MyCleanFunctions, myCleaningParams) #clean the
      doc with the corpusCleaningTools.py
186            DebFile.write ("\nCLEANED TEXT: %s\n\n" % myDOCS[n])
187    if S==1: #stem the words with LSA check
188      print "\n   ...STEMMING THE WORDS...\n"
189      myDOCS = stemm(myDOCS,NameBase)
190      for g,k in myDOCS.iteritems():
191        DebFile.write ("\nSTEMMED TEXT #%d: %s\n\n" % (int(g),k))
192    else:
193      print "\n   !!STEMMING WILL NOT BE PERFORMED!\n"
194
195
196    # The program generates one output file for each LSA spaces in the LSA folder.
197
198    OutFileName = {}
199    DocByDocFile = {}
200    DocByDocFile_Col = {} #additional column file output
201    LSAspaces = {}
202    COS = {}
203
204    print "   ...GENERATING OUTPUT FILES in the directory  %s:\n" % os.path.join(APP_PATH,OUTPUT_PATH)
205    DebFile.write ("\n\n\n*****************************************************************")
206
207    for space in LSAfiles:
208        DocByDocFile[space]=open(os.path.join(APP_PATH,OUTPUT_PATH,'%s_%s_%s.csv' %
      (infilename,stoplist,space)),'w')
209        DocByDocFile_Col[space]=open(os.path.join(APP_PATH,OUTPUT_PATH,'%s_%s_%s_COLUMN.csv' %
      (infilename,stoplist,space)),'w')
210        DocByDocFile[space].write (",")                                    #leaving the first cell of the
      first row empty
211        LSAspaces[space]=lsa(os.path.join(LSA_PATH,space))
212        print "  * FILE: %s_%s_%s.csv" % (infilename,stoplist,space)
213        print "  * FILE: %s_%s_%s_COLUMN.csv" % (infilename,stoplist,space)
214        DebFile.write ("* OUTPUT FILE: %s\n" % os.path.join(APP_PATH,OUTPUT_PATH,'%s_%s_%s.csv' %
      (infilename,stoplist,space)))
215
216    DebFile.write ("\n***************************************************************")
217    DebFile.write ("\n*                        LSA ANALYSIS                         *")
218    DebFile.write ("\n*                   USING SIMILARITY method                   *")
219    DebFile.write ("\n***************************************************************\n\n")
220
221    for counter in range(1, len(myDOCS)+1): #populating line 1 of the output file with DOC numbers
222        for F in LSAfiles:
223            DocByDocFile[F].write ("%d," % counter)
224
225    for key1 in range(1,len(myDOCS)+1): #populating row 1 of the output file with DOC numbers
226        for F in LSAfiles:
227            DocByDocFile[F].write ("\n%d" % key1)
228        DOC1=myDOCS[str(key1)]
229        if DOC1:
230            for key2 in range(1, len(myDOCS)+1):
231                DOC2=myDOCS[str(key2)]
232                if key1<key2: #adding for diagonal
233                    DebFile.write ("\n\nSubject #%d and Subject #%d\n" % (key1, key2))
234                    DebFile.write ("\nSubject #%d para = %s \n\nSubject #%d para = %s \n" %(key1, DOC1,
      key2, DOC2))
235                    if DOC2:
236                        for P in LSAfiles:
237                            COS[P] = LSAspaces[P].Similarity(DOC1, DOC2)
238                            DocByDocFile[P].write (",%f" % COS[P])
239                            DebFile.write ("\nCosine from %s = %f" % (P, COS[P]))
```

```python
240                     DocByDocFile_Col[P].write("%d_%d,%f\n" %(int(key1), int(key2), COS[P]))
241                 else:
242                     for D in LSAfiles:
243                         DocByDocFile[D].write (",")
244                         DocByDocFile_Col[D].write("%d_%d,,\n" %(int(key1), int(key2)))
245                         DebFile.write ("\nCosine from %s = N/A" % D)
246             else:
247                 for K in LSAfiles:
248                     DocByDocFile[K].write (",")
249         else:
250             for key2 in range((key1+1),len(myDOCS)+1):
251                 for U in LSAfiles:
252                     DocByDocFile_Col[U].write("%d_%d,,\n" %(int(key1), int(key2)))
253             for key2 in range(1,len(myDOCS)+1):
254                 for U in LSAfiles:
255                     DocByDocFile[U].write (",")
256                     DebFile.write ("\nCosine from %s = N/A" % U)
257
258     for C in LSAfiles:
259         DocByDocFile[C].close()
260         DocByDocFile_Col[C].close()
261
262     DebFile.close()
263
264     print "\n\nDONE!\n"
265
266
267
268 if __name__ == "__main__":
269     main(sys.argv[1:])
```

# Bibliography

Dennis, S., & Stone, B. (2011). *SEMMOD: Semantic models package* (Version 1.5): Ohio State University.

Landauer, T. K., & Dumais, S. (1997). A Solution to Plato's Problem: The Latent Semantic Analysis Theory of Acquisition, Induction, and Representation of Knowledge. *Prychological Review, 104(2)*, 211-240.

Landauer, T. K., Foltz, P. W., & Laham, D. (1998). Introduction to Latent Semantic Analysis. *Discourse Processes, 25(2-3)*, 259-284.

Stone, B., & Dennis, S. (2012). *Wikipedia Subcorpora Tool* (Version 1.0): Ohio State University.

# List of symbols/abbreviations/acronyms/initialisms

| ARP | Applied Research Project |
|---|---|
| ASSN | Analysis of Semantic and Social Networks |
| CSV | Comma Separated Values |
| DND | Department of National Defence |
| DRDC | Defence Research & Development Canada |
| DRDKIM | Director Research and Development Knowledge and Information Management |
| GUI | Graphical User Interface |
| LSA | Latent Semantic Analysis |
| NLTK | Natural Language Toolkit |
| R&D | Research & Development |
| OSU | Ohio State University |
| SEMMOD | Semantic Models |
| TIF | Technology Investment Fund |
| TN | Technical Note |
| UTF | Unicode Transformation Formats |

This page intentionally left blank.

| | DOCUMENT CONTROL DATA | | |
|---|---|---|---|
| | (Security classification of title, body of abstract and indexing annotation must be entered when the overall document is classified) | | |

**DOCUMENT CONTROL DATA**

(Security classification of title, body of abstract and indexing annotation must be entered when the overall document is classified)

1. ORIGINATOR (The name and address of the organization preparing the document. Organizations for whom the document was prepared, e.g. Centre sponsoring a contractor's report, or tasking agency, are entered in section 8.)

   Defence R&D Canada – Toronto
   1133 Sheppard Avenue West
   P.O. Box 2000
   Toronto, Ontario
   M3K 2C9

2. SECURITY CLASSIFICATION (Overall security classification of the document including special warning terms if applicable.)

   UNCLASSIFIED
   (NON-CONTROLLED GOODS)
   DMC A
   REVIEW: GCEC JUNE 2010

3. TITLE (The complete document title as indicated on the title page. Its classification should be indicated by the appropriate abbreviation (S, C or U) in parentheses after the title.)

   Latent Semantic Analysis (LSA) tools:

4. AUTHORS (last name, followed by initials – ranks, titles, etc. not to be used)

   Derbentseva, N.; Kwantes, P.; Terhaar, P.

5. DATE OF PUBLICATION (Month and year of publication of document.)

   July 2012

6a. NO. OF PAGES (Total containing information, including Annexes, Appendices, etc.)

   35

6b. NO. OF REFS (Total cited in document.)

   4

7. DESCRIPTIVE NOTES (The category of the document, e.g. technical report, technical note or memorandum. If appropriate, enter the type of report, e.g. interim, progress, summary, annual or final. Give the inclusive dates when a specific reporting period is covered.)

   Technical Note

8. SPONSORING ACTIVITY (The name of the department project office or laboratory sponsoring the research and development – include address.)

   Defence R&D Canada – Toronto1133 Sheppard Avenue WestP.O. Box 2000Toronto, Ontario M3K 2C9

9a. PROJECT OR GRANT NO. (If appropriate, the applicable research and development project or grant number under which the document was written. Please specify whether project or grant.)

   15ah

9b. CONTRACT NO. (If appropriate, the applicable number under which the document was written.)

10a. ORIGINATOR'S DOCUMENT NUMBER (The official document number by which the document is identified by the originating activity. This number must be unique to this document.)

   DRDC Toronto TN 2012-079

10b. OTHER DOCUMENT NO(s). (Any other numbers which may be assigned this document either by the originator or by the sponsor.)

11. DOCUMENT AVAILABILITY (Any limitations on further dissemination of the document, other than those imposed by security classification.)

   Unlimited

12. DOCUMENT ANNOUNCEMENT (Any limitation to the bibliographic announcement of this document. This will normally correspond to the Document Availability (11). However, where further distribution (beyond the audience specified in (11) is possible, a wider announcement audience may be selected.))

   Unlimited

13. ABSTRACT (A brief and factual summary of the document. It may also appear elsewhere in the body of the document itself. It is highly desirable that the abstract of classified documents be unclassified. Each paragraph of the abstract shall begin with an indication of the security classification of the information in the paragraph (unless the document itself is unclassified) represented as (S), (C), (R), or (U). It is not necessary to include here abstracts in both official languages unless the text is bilingual.)

NOT REQUIRED.

14. KEYWORDS, DESCRIPTORS or IDENTIFIERS (Technically meaningful terms or short phrases that characterize a document and could be helpful in cataloguing the document. They should be selected so that no security classification is required. Identifiers, such as equipment model designation, trade name, military project code name, geographic location may also be included. If possible keywords should be selected from a published thesaurus, e.g. Thesaurus of Engineering and Scientific Terms (TEST) and that thesaurus identified. If it is not possible to select indexing terms which are Unclassified, the classification of each should be indicated as with the title.)

Latent Semantic Analysis, LSA, semantic models, semantic analysis, python

**Defence R&D Canada**

Canada's Leader in Defence
and National Security
Science and Technology

**R & D pour la défense Canada**

Chef de file au Canada en matière
de science et de technologie pour
la défense et la sécurité nationale

DEFENCE **R&D** DÉFENSE

**www.drdc-rddc.gc.ca**