

gen_bigrams.prl

```
1  #!/usr/local/bin/perl
2
3  #####
4  ## This perl script is used to generate all the bigrams
5  ## within the .stemmed files and .tokenized files (except
6  ## the files for interactive queries which will be generated
7  ## dynamically each time a use make inputs.
8  ##
9  ## @INPUT          @OUTPUT
10 ## cacm.stemmed      >      cacm.stemmed.bigrams
11 ## cacm.stemmed.hist >      cacm.stemmed.hist.bigrams
12 ## cacm.tokenized    >      cacm.tokenized.bigrams
13 ## cacm.tokenized.hist >    cacm.tokenized.hist.bigrams
14 ## query.stemmed     >      query.stemmed.bigrams
15 ## query.stemmed.hist >    query.stemmed.hist.bigrams
16 ## query.tokenized   >      query.tokenized.bigrams
17 ## query.tokenized.hist >  query.tokenized.hist.bigrams
18 ##
19 #####
20
21 use Carp;
22 use FileHandle;
23
24 my $DIR = "/home/pan/test/hw2";
25
26 # Initialize the stopword list for both stemmed and tokenized
27 my $stoplist = "$DIR/common_words";
28 my $stoplist_stemmed = "$DIR/common_words\.stemmed";
29 my %stoplist_hash = ();
30 my %stoplist_stemmed_hash = ();
31
32 # Initialize the tokenized stopword list
33 my $stoplist_fh = new FileHandle $stoplist , "r"
34   or croak "Failed [$stoplist_fh] $stoplist";
35
36 while (defined( $line = <$stoplist_fh> )) {
37   chomp $line;
38   $stoplist_hash{ $line } = 1;
39 }
40
41 # Initialize the stemmed stopword list
42 my $stoplist_stemmed_fh = new FileHandle $stoplist_stemmed , "r"
43   or croak "Failed [$stoplist_fh] $stoplist_stemmed";
44
45 while (defined( $line = <$stoplist_stemmed_fh> )) {
46   chomp $line;
47   $stoplist_stemmed_hash{ $line } = 1;
48 }
49
50
51 &gen_bigrams("cacm");
52 &gen_bigrams("query");
53
54 exit(0);
55
56
57 #####
58 ## GEN_BIGRAMS
59 ##
60 ## This function will generate the bigrams term sets
61 #####
62 sub gen_bigrams {
63
64   my $file_name = shift;
65
66   # array to store all the bigrams as value
67   my @bigrams = ();
68   my $last_word = undef; # used to store the last word so we could have pairs
69
70   # my $file_name      = "interactive";
71   my $raw_file         = "$file_name\.raw";
72   my $stemmed_file     = "$file_name\.stemmed";
73   my $tokenized_file   = "$file_name\.tokenized";
74   my $stem_hist_file   = "$file_name\.stemmed\.hist";
75   my $token_hist_file  = "$file_name\.tokenized\.hist";
76   my $stem_hist_com    = "cat $stemmed_file\.bigrams | perl $DIR/make_hist.prl > $stem_hist_file\.bigrams";
77   my $token_hist_com   = "cat $tokenized_file\.bigrams | perl $DIR/make_hist.prl > $token_hist_file\.bigrams";
78
79   my $tokenized_file_fh = new FileHandle $tokenized_file, "r"
80     or croak "Failed $tokenized_file";
81
82   my $stemmed_file_fh = new FileHandle $stemmed_file, "r"
83     or croak "Failed $stemmed_file";
84
85
86   #####
87   # make the stemmed bigrams
88   #####
89   open(RAWFILE, ">$stemmed_file\.bigrams");
```

```

90
91 while (defined( $word = <$stemmed_file_fh> )) {
92     chomp $word;
93
94     # if we encounter a normal word, we check if it eligible for bigram
95     if ($last_word =~ /^[a-zA-Z]/ and $word =~ /^[a-zA-Z]/) { # start of query tokens
96         # if both words are not in the stoplist, we added to the bigrams array
97         if (! exists $stoplist_stemmed_hash{ $last_word } and ! exists $stoplist_stemmed_hash{ $word }) {
98             push (@bigrams, "$last_word+$word");
99         }
100     }
101
102     # if next doc/query, pop all the elements in the array into file
103     if ($word =~ /\^[A-Z]/ and $#bigrams != -1) {
104         # after we scanned each doc/qry, all the bigrams are in the array now
105         # we append all the bigrams to the end of doc/qry
106         foreach $bigram (@bigrams) {
107             print RAWFILE "$bigram\n";
108         }
109         # clean up the array
110         @bigrams = ();
111     }
112
113     # we shift last_word by 1 and write the current word into the output file
114     $last_word = $word;
115     print RAWFILE "$word\n";
116 }
117
118 # do it once again to make sure we have everything write into the file
119 if ($#bigrams != -1) {
120     foreach $bigram (@bigrams) {
121         print RAWFILE "$bigram\n";
122     }
123 }
124
125 # close the bigram file after we created it
126 close(RAWFILE);
127
128 # generate hist file
129 system ("${stem_hist_com}") and die "Failed $DIR/make_hist.prl: ${!}\n";
130
131 # clean up the array
132 @bigrams = ();
133
134 #####
135 # make the tokenized bigrams
136 #####
137 open(RAWFILE, ">${tokenized_file}.bigrams");
138
139 while (defined( $word = <$tokenized_file_fh> )) {
140     chomp $word;
141
142     # if we encounter a normal word, we check if it eligible for bigram
143     if ($last_word =~ /^[a-zA-Z]/ and $word =~ /^[a-zA-Z]/) { # start of query tokens
144         # if both words are not in the stoplist, we added to the bigrams array
145         if (! exists $stoplist_hash{ $last_word } and ! exists $stoplist_hash{ $word }) {
146             push (@bigrams, "$last_word+$word");
147         }
148     }
149
150     # if next doc/query, pop all the elements in the array into file
151     if ($word =~ /\^[A-Z]/ and $#bigrams != -1) {
152         # after we scanned each doc/qry, all the bigrams are in the array now
153         # we append all the bigrams to the end of doc/qry
154         foreach $bigram (@bigrams) {
155             print RAWFILE "$bigram\n";
156         }
157         # clean up the array
158         @bigrams = ();
159     }
160
161     # we shift last_word by 1 and write the current word into the output file
162     $last_word = $word;
163     print RAWFILE "$word\n";
164 }
165
166 # do it once again to make sure we have everything write into the file
167 if ($#bigrams != -1) {
168     foreach $bigram (@bigrams) {
169         print RAWFILE "$bigram\n";
170     }
171 }
172
173 # close the bigram file after we created it
174 close(RAWFILE);
175
176 # generate hist file
177 system ("${tokn_hist_com}") and die "Failed $DIR/make_hist.prl: ${!}\n";
178
179 # clean up the array
180 @bigrams = ();

```

