

# CS231 Project 8

## Summary

The goal of this project was to find the most common words in a text document of reddit comments and find trends related to how many times a set of words appeared over 8 years of reddit comments. To find the most common words, counts files that include each word and how often it appeared were found in Project 7 using a binary search tree and were read into a binary search tree again, then put into a priority queue. The priority queue removes items in order of how many times the word appears and prints the list to the terminal. In an extension, it reads the word-value pairs directly into a priority queue. Then to find trends, it reads the word-value pairs into a binary search tree and finds each word in a list and returns its frequency (how many times it appears divided by the total number of words in the document) and prints a list of how many times it appears to the terminal. Then, a graph is made with the results. In extensions, it prints the results to files rather than printing to the terminal. The most common words found are pronouns articles and very simple verbs like is and have. Cat and dog are the most common pet words over all 8 years tested.

## Tasks

1) The first task was to create a FindCommonWords task that will print the words in a word count file (from project 7) in order from the highest number they occur to the least number of times they occur. It has fields for a word counter object (from project 7) and a priority queue heap that uses a comparator that takes in key value pairs that have a string (the word) and an integer (the count) and returns the difference between the counts to find which word occurs more times.

In findWords, I used the readWordCountFile of WordCounter from Project 7 to put all of the word-count pairs into a binary search tree. Then using the getPairs method, I made an arrayList of the word-count pairs. Then I looped over every pair in the arrayList and add them to the priority queue heap. Then it loops over the number of items in the heap and prints out the word-count pair returned when removed is called. This will print every word in the word count file in order from the word that occurs most to the word that occurs least. Here is an image of the ten most common words in the 2008 reddit comments file:

```
Key: the, Value: 715631
Key: to, Value: 433159
Key: a, Value: 383959
Key: i, Value: 350554
Key: of, Value: 331367
Key: and, Value: 320678
Key: that, Value: 290574
Key: it, Value: 277506
Key: is, Value: 259969
Key: you, Value: 246681
```

"the" is by far the most common word, almost twice more common than even the second most common word.

"deleted" is the first non pronoun, article, auxiliary or linking verb, word that appears in the list, the 27th occurring 74,414 times.

2) The second task was to create a findTrends class to "determine the frequency of a specific set of words in a set of words". It has the same fields as the findCommonWords class.

I made a findTrends method that starts by clearing the word counter object in case more than one list of words is analyzed with the same FindTrends object. Then it calls the readWordCountFile method of the WordCounter object. It makes an arrayList of Doubles called frequency and puts 0.0 in each location up to the size of the words arrayList passed in to the method. Then for each word in the arrayList, it calculates the frequency by calling the getFrequency method of the Word Counter object and then sets the index of the word in the frequency array list to the frequency of the corresponding word in the words array list. At the end it returns the frequency array list.

Then I used a secondary helper method called findTrend with an array of strings (for the args array) and an array list of words. It starts by making an arrayList frequency of word by Year with empty strings the length of the number of words in the words arrayList. Then looping over the 1st index of the args array up to and including the 2nd index of the args array, it calls findTrends from above on the current year file (args[0]i.txt).

Then it loops over that arrayList and sets the corresponding index in the frequency of word by year to its current value plus the index of the findTrends arrayList. At the end, I print the ArrayList.

I put in a list of animals (cat, dog, hamster, fish, rabbit, frog, ferret, hedgehog, turtle, gecko, lizard, alpaca). Here is the terminal output with the frequencies from the reddit comments in 2008-2015:

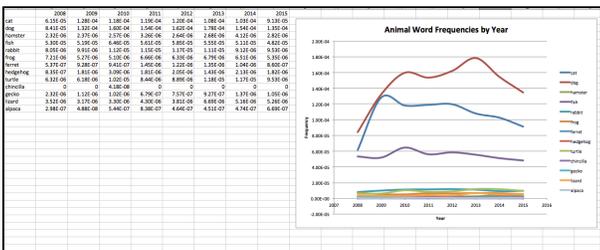
```

kaylas-mbp:Project8 kaylafreemans java FindCommonWords counts_reddit_comments_2008 2015 cat dog hamster fish ra
bbat frog ferret hedgehog turtle chincilla gecko lizard alpaca
cat,6,151911804952444E-5,1.2018046245149004E-4,1.180697447537388E-4,1.1891344486613104E-4,1.196664902976222E-4,
1.8789311480197898E-4,1.40255984153580734E-4,0.1216657846724159E-5
dog,8,485228338161769E-5
hamster,3,3248583914864468E-6
fish,5,299466661436747E-5
rabbit,8,847559047176163E-6
frog,7,212997380280746E-6
ferret,5,365839364784108E-7
hedgehog,8,345616789664169E-7
turtle,6,318924148745727E-6
chincilla,0,0
gecko,2,3248583914864468E-6
lizard,3,517081361358471E-6
alpaca,2,988574248898E-7
cat,6,151911804952444E-5,1.2018046245149004E-4,1.180697447537388E-4,1.1891344486613104E-4,1.196664902976222E-4,
1.8789311480197898E-4,1.40255984153580734E-4,0.1216657846724159E-5
dog,8,485228338161769E-5,1.3221193720001334E-4,1.5969100177892724E-4,1.5352579156331634E-4,1.6204684560817018E-4,
1.703856199095687E-4,1.544082231915196E-4,1.3463013732536818E-4
hamster,3,3248583914864468E-6,2.367982124889156E-6,2.5712788755883193E-6,3.2597631945027354E-6,2.638085517890018
3E-6,2.6884186673674834E-6,4.1213423619466993E-6,2.819229140807568E-6
fish,5,299466661436747E-5,5.187414448855767E-5,6.455371298519493E-5,6.072454394328384E-5,5.847756231322074E-5,5.5
5487174788204915E-5,5.114111963848895E-5,4.816581318524885E-5
rabbit,8,847559047176163E-6,9.911013017578549E-6,1.1225823793872996E-5,1.1454445669572111E-5,1.1724824523955637E
-5,1.1072383653985395E-5,9.119862122698884E-6,9.532817177075759E-6
frog,7,212997380280746E-6,5.272854216248886E-6,5.108746751414877E-6,6.655349855443084E-6,6.326519899384395E-6,6.
788723914547552E-6,6.5136158081927718E-6,5.351757011691654E-6
ferret,5,365839364784108E-7,9.276317682659725E-7,9.48711491839629E-7,1.4487836428012157E-6,1.2213358879128454E-6,
1.352734545592997E-6,1.0421785383884148E-6,8.481838954584444E-7
hedgehog,8,345616789664169E-7,1.8864487963074202E-6,3.093895578530335E-6,1.81097955250515196E-6,2.85184429169223
62E-6,1.4278865798125846E-6,2.1317288079036167E-6,1.8157747803953824E-6
turtle,6,318924148745727E-6,6.17869748514882396E-6,1.892238026318635E-5,8.443692163538335E-5,8.891325263999691E-6,
1.1823982986518244E-5,1.16534981653977E-5,9.532817177075759E-6
chincilla,0,0,0,4,1.88939968176129E-8,0,0,0,0,0,0,0,0,0
gecko,2,3248583914864468E-6,1.122920265716487E-6,1.8243382982431515E-6,6.79173321888699E-7,7.572282585054682E-7,
9.2697374798625E-7,1.3737809791356641E-6,1.85137984439432E-6
lizard,3,517081361358471E-6,3.1734770745941164E-6,3.3829425885391416E-6,4.3810764371911894E-6,3.8185679828582E
-6,6.68852124754318E-6,5.16352888925427E-6,5.2561899221971E-6
alpaca,2,988574248898E-7,4.882272422452487E-8,5.435221948228967E-7,8.375780430319528E-7,4.641876374865773E-7,4.
5091155197638E-7,4.737175128674784E-7,6.689696264614568E-7
kaylas-mbp:Project8 kaylafreemans java FindCommonWords counts_reddit_comments_2008 2015 cat dog hamster fish ra
bbat frog ferret hedgehog turtle chincilla gecko lizard alpaca

```

It shows the animals in just 2008, then in 2008-2015.

I copied and pasted this into an excel document, added a header, then made a graph of the results:



Dog, cat, and fish are by far the most common animal names uses throughout all 8 years. The three words are closest in 2008 then vary but generally grow apart until 2015. Dog and cat become generally more frequent between 2008 and 2015, while fish remains constant. Cat almost catches up to dog in 2009, but not quite. The rest of the words are all much less common and remain constant in their frequentness over then 8 years.

## Extensions

1) The first extension that I completed was to move the findCommonWords output to as text file rather than printing to the terminal so it is easier to look through and understand. To do this, I made a new method called writeFindWordsFile. It starts by clearing the WordCount object, then calling readWordCountFile of the WordCount object. It makes an ArrayList of the pairs then adds each to the priority queue heap. Then I made a BufferedWriter object and in a try loop, I made a new bufferedWriter object and looped over the number of items in the priority queue heap, and writes the key value pair returned by removing from the priority queue heap). Then it catches exceptions and will close the file when necessary.

```
Key: the, Value: 715631
Key: to, Value: 433159
Key: a, Value: 383959
Key: i, Value: 350554
Key: of, Value: 331367
Key: and, Value: 320678
Key: that, Value: 290574
Key: it, Value: 277506
Key: is, Value: 259969
Key: you, Value: 246681
Key: in, Value: 219746
Key: s, Value: 172222
Key: t, Value: 159511
Key: for, Value: 148217
Key: not, Value: 118092
Key: are, Value: 115588
Key: they, Value: 110612
Key: be, Value: 110401
Key: this, Value: 109081
Key: on, Value: 105931
Key: have, Value: 104203
Key: but, Value: 98811
Key: with, Value: 96434
Key: as, Value: 89151
Key: if, Value: 85114
Key: was, Value: 83278
Key: deleted, Value: 74414
Key: or, Value: 72585
Key: can, Value: 69404
Key: what, Value: 68586
```

Here is the text file made using this method:

(I tried to upload the file but its too big)

2) The second extension that I completed was to try to make the findWords process faster. I made another method that would create a text file, called findWords that requires an integer as a parameter to distinguish it from the first. It starts by initiating a String object and a BufferedReader object. In a try loop, it makes a BufferedReader object, then uses it to read the next line of the file (to eliminate the header). Then it reads each line of the word counts file and parses it, making the two pieces into a KeyValuePair object, then adding each KeyValuePair to the priority heap. Then it checks for exceptions and closes the file. In findWords, it prints as it removes, in writeFindWordsNewFile, it makes a text file using this new method, so then it writes the removes to a text file just like the previous extension did. I made sure that the file created was the same as above:

```
Key: the, Value: 715631
Key: to, Value: 433159
Key: a, Value: 383959
Key: i, Value: 350554
Key: of, Value: 331367
Key: and, Value: 320678
Key: that, Value: 290574
Key: it, Value: 277506
Key: is, Value: 259969
Key: you, Value: 246681
Key: in, Value: 219746
Key: s, Value: 172222
Key: t, Value: 159511
Key: for, Value: 148217
Key: not, Value: 118092
Key: are, Value: 115588
Key: they, Value: 110612
Key: be, Value: 110401
Key: this, Value: 109081
Key: on, Value: 105931
Key: have, Value: 104203
Key: but, Value: 98811
Key: with, Value: 96434
Key: as, Value: 89151
Key: if, Value: 85114
Key: was, Value: 83278
Key: deleted, Value: 74414
Key: or, Value: 72585
Key: can, Value: 69404
Key: what, Value: 68586
```

Then, to see which worked best, I ran them on the same file (2008) and kept track of how long it took to do each:

```
kaylas-mbp:Project8 kaylafreeman$ java FindCommonWords counts_reddit_comments_20
08.txt counts_reddit_comments_2008_times.txt
Directly Into PQHeap: 39268
Make Binary Tree: 38652
kaylas-mbp:Project8 kaylafreeman$ java FindCommonWords counts_reddit_comments_20
08.txt counts_reddit_comments_2008_times.txt
Directly Into PQHeap: 38254
Make Binary Tree: 38343
kaylas-mbp:Project8 kaylafreeman$ java FindCommonWords counts_reddit_comments_20
08.txt counts_reddit_comments_2008_times.txt
Directly Into PQHeap: 38889
Make Binary Tree: 39104
```

The two options actually take approximately the same amount of time, so not much was saved by adding KeyValuePairs directly into the priority queue heap.

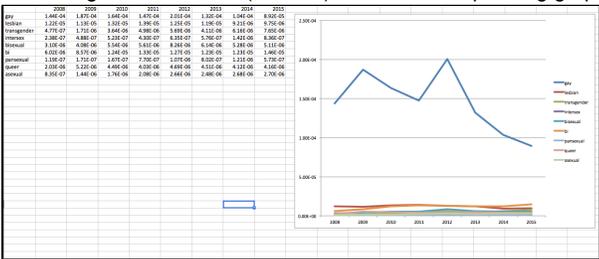
3) The third extension that I completed was to write the results of findTrend to a CSV file so it can be opened as a CSV file. I edited findTrend so that at the end I initialized a BufferedWriter object to null. In a try loop, I made a new bufferedWriter, I wrote a header by first writing the numbers in the range of the first index in the args array up to and including the second index of the args array, separated by commas. This makes a header with the list of years that the frequency information was collected at. Then I write every item in the words list followed by the corresponding string in the frequency of words by year. This makes a full csv file with a header line that can be turn into a graph. It looks like this:

```

| , 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015
cat:
6.151911804952444E-5, 1.2818406245149004E-4, 1.1806974447537308E-4, 1.1891344486613104E-4, 1.1
9664902976222E-4, 1.0789311402197898E-4, 1.0255984153580734E-4, 9.126657046724159E-5
dog,
8.485228338161769E-5, 1.3221193720001334E-4, 1.5969100177892724E-4, 1.5352579156331634E-4, 1.6
204684560817018E-4, 1.783856199895687E-4, 1.544082231915196E-4, 1.346301373253681E-4
hamster,
2.3248503914064468E-6, 2.367902124889456E-6, 2.5712780755083193E-6, 3.2597631945027354E-6, 2.6
308855178900183E-6, 2.6804186673674834E-6, 4.121342361946993E-6, 2.819229140087568E-6
fish,
5.299466661436747E-5, 5.187414448855767E-5, 6.455371298511943E-5, 5.6072454394328304E-5, 5.847
756231322874E-5, 5.5487171478682015E-5, 5.114119638400056E-5, 4.8165813105224885E-5
rabbit,
8.847559847176163E-6, 9.911813017578549E-6, 1.1225823793072906E-5, 1.145444566957211E-5, 1.17
24824523955637E-5, 1.1072383653985305E-5, 9.119062122668804E-6, 9.532817177075759E-6
frog,
7.212997368209746E-6, 5.272854216248686E-6, 5.100746751414877E-6, 6.655349855443084E-6, 6.3265
10899384395E-6, 6.788723914547552E-6, 6.513615801927718E-6, 5.351757011691654E-6
ferret,
5.365039364784108E-7, 9.276317602659725E-7, 9.40711491039629E-7, 1.4487836420012157E-6, 1.2213
358879120454E-6, 1.3527346545592907E-6, 1.8421785283084348E-6, 8.601038054584444E-7
hedgehog,
8.345616789664169E-7, 1.8064407963074202E-6, 3.0938955705303353E-6, 1.8109795525015196E-6, 2.0
5184291926232E-6, 1.4278865798125846E-6, 2.1317288079036167E-6, 1.8157747003953824E-6
turtle,
2.3248503914064468E-6, 1.122922657164072E-6, 1.0243302902431515E-6, 6.791173321800699E-7, 7.57
2282505854682E-7, 9.268737447986251E-7, 1.3737807873156641E-6, 1.051237984439432E-6
lizard,
3.517081361358471E-6, 3.1734770745941164E-6, 3.3029425685391416E-6, 4.3010764371911094E-6, 3.8
10567970285582E-6, 6.68852134754316E-6, 5.163520890255427E-6, 5.25618992219716E-6
alpaca,
2.98857742488086E-7, 4.802272422452487E-8, 5.435221948228967E-7, 8.375780438319528E-7, 4.64107
6374065773E-7, 4.509115515197636E-7, 4.737175128674704E-7, 6.689696264614568E-7

```

4) The last extension that I completed was to test FindTrends on another set of words. I decided to do a list of words related to sexuality to see how their frequency has changed in the last 8 years. I used gay, lesbian, transgender, intersex, bisexual, bi, pansexual, queer and asexual. Here is an image of the csv file (in excel) and the corresponding graph:



Gay is by far the most common word used, although its frequency has decreased since 2012. I wonder if this is because there has been a lot of talk of eliminating "gay" as an insult or because people now use other words to discuss sexuality. Lesbian is the only other word to generally decrease in frequency over the 8 years, and only very slightly. The rest of the words increase in frequency slightly, possibly because there is more knowledge and use of these words in modern culture.

5) The fifth extension I completed was to make a findWords text file for each of the reddit comment counts file (for each year from 2008-2015). Interestingly, they all had the same top 10 words, with the and i always being the top two, and the rest occasionally switching positions.

**Help**

Prof: Stephanie

Students: Jos