Implementind Levenshtein Distance: (Extra Credit)

The corpus contains a lot of files containing the given keywords and not just 23.

 The Levenshtein distance gives the edit distance between every word in the file and given keywords. If this distance is 0(exact match) or 1(differs only by either 1 insertion, deletion or replacement) then, the total number of such other words found is calculated and placed in the output file and all these words are also stored in the log file.

 Clearly the ouput file shows that there are a large number of words that differ only by 1 edit distance from the given keywords and running the program on these new words would yield more files which would include irrelevant files because eg. words with levenshtein distance 1 from ufo are 2545 which include uf, uno, ifo, wfo ... which don't really make sense as ufo. So implementing levenshtein distance on the keywords and all the words of the files wouldn't be beneficial in increasing the number of relevant files for this program.

```
Keyword(s) used: UFO, flying disc, disc, saucer, extraterrestrial craft, flying saucer, space ship, space craft, air ship, phantom, space probe
No of files processed: 2067
No of files containing keyword(s): 120

No of occurrences of each keyword:
--------------------------------
        UFO: 158
        flying disc: 15
        disc: 113
        saucer: 120
        extraterrestrial craft: 0
        flying saucer: 30
        space ship: 4
        space craft: 3
        air ship: 0
        phantom: 1
        space probe: 0
No of files containing each keyword:
--------------------------------
        UFO: 65
        flying disc: 4
        disc: 55
        saucer: 8
        extraterrestrial craft: 0
        flying saucer: 6
        space ship: 3
        space craft: 3
        air ship: 0
        phantom: 1
        space probe: 0
No of matches with Levenshtein Distance for each keyword:
--------------------------------
        UFO: 2545
        flying disc: 1
        disc: 523
        saucer: 169
        extraterrestrial craft: 0
        flying saucer: 0
        space ship: 1
        space craft: 6
        air ship: 21
        phantom: 1
```