

Web Crawler

The topic I covered for my project was the Japanese magazine publisher Nikoli, who are better known as the creators of Sudoku, along with many other pencil and paper puzzles. I created my knowledge base with no modifications, purely through the execution of the code. I began with a list of three websites: a page on Nikoli's website listing their puzzles and instructions on how to solve them, the Wikipedia page for Nikoli, and the page for Nikoli on the U.S. based parent publisher Hachette Book Group. After lots of testing, I found that these three sites tended to give the best output and most suitable important terms.

I initialize three variables: `num_pages_to_scrape`, `linkspersdomain`, and `linksperspage`. The first controls the total number of pages scraped/visited by the webcrawler, the second controls the number of pages allowed per each base-level web domain, and links per page controls the number of links added to the list of links by each webpage.

The program keeps track of a list of URLs aptly titled "urls". The list of 3 base URLs are added to this list, and during the crawling process, links found on each page are appended to this list according to a list of criteria which will be explained later. For each URL, it loads it with BeautifulSoup, removes unwanted tags like headers and footers, unwraps formatting like italics and bold text, and extracts the text. After doing this, it removes foreign characters (this dataset in particular has many pages with Japanese text, so they must be regex'ed away), saves the modified text to a file, and adds the URL to a list of already traversed URLs.

Next, it gathers all the links on the page by looking up all the href tags, filtering out any that do not begin with `https://`, and ignoring any additional tags if more than the number in `linksperspage` have been looked up. The program then finds a rough approximation of the domain, grabbing the fifth to the fifteenth characters, and adds that to a dictionary, incrementing the value by 1 if it already is present. This list is used to fulfill the domain check which also uses `linkspersdomain`. The program also filters out keywords in URLs corresponding to domains such as Wikimedia and Google, which

were determined by which sites were giving me the most trouble in testing. If the link passes all of these checks, it is added to the end of the URL list, incrementing the number used to check against linkssopage (defined as the variable “linkssopar”).

Once it loops through this process enough times, outputting as many files as determined in num_pages_to_scrape, the files are then cleaned for NLP parsing. All the Nikoli text is lowercased, with all lines shorter than 15 characters being removed from the dataset (hopefully to get as many untagged titles and footers as possible). I forewent removing stopwords, since the TF-IDF process I use shows that they are common to all documents anyways. All the files are concatenated into one long nikoli_corpus string.

This nikoli_corpus string is compared to the text from three other unrelated documents: two are chapters of textbooks used in the homework, and the other is the first 15000 characters of Moby Dick, which is one of the default texts in NLTK. The TF-IDF is calculated for the Nikoli text against the three other documents, with the terms most specific to the Nikoli text being sorted to the top. This returns a large list of words, of which I picked the following for my important terms:

[“nikoli”, “puzzle”, “sudoku”, “rules”, “book”, “magazine”, “japan”, “solve”, “cells”, block”]

These terms are then looked up in the unmodified corpus for sentences containing them—sentences are segmented using the NLTK sentence tokenizer. This tokenizer isn’t perfect, however, and often it will mistag long unbroken strings as single sentences, so any sentence longer than 400 characters is thrown out, line returns are removed, and the sentence is added to a list, which is then the value for each word’s key in the knowledge base dictionary. This knowledge base is pickled and saved to a file to be used by the chatbot.

Knowledge base screenshots

```

1 {
2   'nikoli': [
3     'nikoli book shop page "Puzzle Communication Nikoli Vol.186" is now on sale.',
4     'All our magazines and books are available at the "nikoli book shop" page.',
5     'You can even access Nikoli puzzles online at nikoli.com.',
6     'Nikoli Video / SNS We will deliver the latest information on nikoli.',
7     'Let's Begin nikoli Puzzles!',
8     'Selected 200 good Sudoku puzzles from Nikoli, the original home of Sudoku.nikoli Place light bulbs (circles) according to the following rules.',
9     'All our magazines and books are available at the "nikoli book shop" page.',
10    'You can order Nikoli books from Nikoli formal web site, nikoli.co.jp .',
11    'We have another exciting puzzle site, nikoli.com, where anyone, anywhere can access and play our unique challenging puzzles.',
12    'Let's Begin nikoli Puzzles!',
13    'Learn more about this authornikoli Place balloons (white circles) and iron balls (black circles) according to the following rules.',
14    'How YouTube works Test new features NFL Sunday Ticket © 2024 Google LLC - YouTubeNikoli Place numbers in squares (parts of snakes) so that the figures fr
15    'With Nikoli logic puzzles you need just paper, pencil, and an active mind, no expensive programs or training.... nikoli Puzzles page nikoli book shop "Puzzl
16    'Learn more about this authornikoli Divide the grid into blocks of five cells.',
17    'nikoli book shop page "Puzzle Communication Nikoli Vol.186" is now on sale.',
18    'We are confident it will be Nikoli puzzles.nikoli A rectangle, bordered by bold lines, is called a "room".',
19    'Nikoli Video / SNS We will deliver the latest information on nikoli.',
20    'All Rights Reservednikoli Fill in the cells under the following rules.',
21    'Effective Date: March 31, 2008 Last Updated: November 1, 2023nikoli Fill in all cells, each with a - or a |.',
22    'With Nikoli logic puzzles you need just paper, pencil, and an active mind, no expensive programs or training.... nikoli Puzzles page nikoli book shop "Puzzl
23  'puzzle': [
24    'The joy of solving puzzles is more than a passing fad, it enriches lives.',
25    'Creating a puzzle does not necessarily require any special skill either.',
26    'You also have a chance to see Nikoli puzzles through video games and books in your country.',
27    'What makes solving puzzles fun?',
28    'Nikoli - The leading company of the puzzle-world.',
29    'An example of a -dependent genre of puzzle would be the , which relies on a specific language and .',
30    'Don't settle when it comes to your Sudoku: Choose puzzles hand-crafted by Nikoli, the Japanese company that invented the game.',
31    'These puzzles with numbers and simple symbols do not depend on a particular language for solving.',
32    '"Nikoli Puzzle Book Reprint Edition Slitherlink 2" The puzzles in this volume are all the puzzles of "Pencil Puzzle Book Slitherlink Vol.11".',
33    'A unique aspect of Nikoli puzzles is that a puzzle that is challenging to a beginner and yet solvable, can also give the expert solver pleasure in solving
34    'Our carefully reviewed puzzles will not cause solvers unnecessary stress.',
35    '"Nikoli Puzzle Book Reprint Edition Kakuro 2" The puzzles in this volume are all the puzzles of "Pencil Puzzle Book Kakuro Vol.11".',
36    'The magazine has invented several new genres of puzzles, and introduced several new games to Some of the popular Nikoli puzzles, along with their names: t
37    'They hold the trademark on sudoku in Japan, and supply 100 percent of the puzzles published in Japanese newspapers and magazines.',
38    'It is said that more than a hundred million people are solving sudoku puzzles every day.',
39    'You can solve these puzzles without knowing Japanese or mathematical techniques.',
40    'The bond between puzzle creator and solver.',
41    'We just wish to share the joy and excitement of puzzles with our puzzle creators and solvers.',
42    'The Nikoli pages on Wikipedia are very useful to understand the rules of Nikoli puzzles for English speakers.',
43  ],
44  'rules': [
45    'Fill in all empty cells with numbers under the following rules.',
46    'The Nikoli pages on Wikipedia are very useful to understand the rules of Nikoli puzzles for English speakers.',
47    'Fill in cells under the following rules.',
48    'Akari has very simple rules and guide you to a logical complex world YosenabeNikoli has been creating puzzles for 40 years and is surely
49    'Learn more about this authornikoli Place balloons (white circles) and iron balls (black circles) according to the following rules.',
50    'Here on our homepage you get the basic rules at each puzzle.',
51    'Paint some of the cells black under the following rules.',
52    'All Rights Reservednikoli Fill in the cells under the following rules.',
53    'Selected 200 good Sudoku puzzles from Nikoli, the original home of Sudoku.nikoli Place light bulbs (circles) according to the following
54  ],
55  'book': [
56    'Then there is the new generation, young people becoming Nikoli fans by discovering the puzzles on the bookshelves of their parents.',
57    'With Nikoli logic puzzles you need just paper, pencil, and an active mind, no expensive programs or training.... nikoli Puzzles page nikoli
58    'Each puzzle is handcrafted (not computer generated, like most generic sudoku books) by the experts at Nikoli Publishing,\xa0the Japanese
59    'nikoli book shop page "Puzzle Communication Nikoli Vol.186" is now on sale.',
60    'Novel Suspects Logo Product image pagination Open the full-size image Nikoli Publishing Formats and Prices Also available from: The ultis
61    'All our magazines and books are available at the "nikoli book shop" page.',
62    'With Nikoli logic puzzles you need just paper, pencil, and an active mind, no expensive programs or training.... nikoli Puzzles page nikoli
63    'You also have a chance to see Nikoli puzzles through video games and books in your country.',
64    'All our magazines and books are available at the "nikoli book shop" page.',
65    'nikoli book shop page "Puzzle Communication Nikoli Vol.186" is now on sale.',
66    'You can order Nikoli books from Nikoli formal web site, nikoli.co.jp .',
67    'We publish a quarterly magazine called Puzzle Communication Nikoli along with many other books.',
68    'And unlike Workman's previous small-sized sudoku books, The Biggest Bestest Book of Sudoku\u200b will come in a larger 8"x10" format and
69  ],
70  'magazine': [
71    '"Puzzle Communication Nikoli Vol.186" Nikoli is the most advanced pencil puzzle magazine in Japan.',
72    'We started the puzzle magazine Puzzle Communication Nikoli in 1980, and Nikoli has always specialized in logic puzzles.',
73    'All our magazines and books are available at the "nikoli book shop" page.',
74    'They hold the trademark on sudoku in Japan, and supply 100 percent of the puzzles published in Japanese newspapers and magazines.',
75    '"Puzzle Communication Nikoli Vol.186" Nikoli is the most advanced pencil puzzle magazine in Japan.',
76    'We publish a quarterly magazine called Puzzle Communication Nikoli along with many other books.',
77    'The magazine has invented several new genres of puzzles, and introduced several new games to Some of the popular Nikoli puzzles, along wi
78    'We also provide puzzles to major Japanese newspapers and magazines.',
79    'All our magazines and books are available at the "nikoli book shop" page.',
80    'They hold the trademark on sudoku in Japan, and supply 100 percent of the puzzles published in Japanese newspapers and magazines.'],
81  'Japan': [
82    'They hold the trademark on sudoku in Japan, and supply 100 percent of the puzzles published in Japanese newspapers and magazines.',
83    'Don't settle when it comes to your Sudoku: Choose puzzles hand-crafted by Nikoli, the Japanese company that invented the game.',
84    .....
85  ],
86  .....
87  },
88  .....
89  },
90  .....
91  },
92  .....
93  },
94  .....
95  },
96  .....
97  },
98  .....
99  },
100  .....
101  },
102  .....
103  },
104  .....
105  },
106  .....
107  },
108  .....
109  },
110  .....
111  },
112  .....
113  },
114  .....
115  },
116  .....
117  },
118  .....
119  },
120  .....
121  },
122  .....
123  },
124  .....
125  },
126  .....
127  },
128  .....
129  },
130  .....
131  },
132  .....
133  },
134  .....
135  },
136  .....
137  },
138  .....
139  },
140  .....
141  },
142  .....
143  },
144  .....
145  },
146  .....
147  },
148  .....
149  },
150  .....
151  },
152  .....
153  },
154  .....
155  },
156  .....
157  },
158  .....
159  },
160  .....
161  },
162  .....
163  },
164  .....
165  },
166  .....
167  },
168  .....
169  },
170  .....
171  },
172  .....
173  },
174  .....
175  },
176  .....
177  },
178  .....
179  },
180  .....
181  },
182  .....
183  },
184  .....
185  },
186  .....
187  },
188  .....
189  },
190  .....
191  },
192  .....
193  },
194  .....
195  },
196  .....
197  },
198  .....
199  },
200  .....
201  },
202  .....
203  },
204  .....
205  },
206  .....
207  },
208  .....
209  },
210  .....
211  },
212  .....
213  },
214  .....
215  },
216  .....
217  },
218  .....
219  },
220  .....
221  },
222  .....
223  },
224  .....
225  },
226  .....
227  },
228  .....
229  },
230  .....
231  },
232  .....
233  },
234  .....
235  },
236  .....
237  },
238  .....
239  },
240  .....
241  },
242  .....
243  },
244  .....
245  },
246  .....
247  },
248  .....
249  },
250  .....
251  },
252  .....
253  },
254  .....
255  },
256  .....
257  },
258  .....
259  },
260  .....
261  },
262  .....
263  },
264  .....
265  },
266  .....
267  },
268  .....
269  },
270  .....
271  },
272  .....
273  },
274  .....
275  },
276  .....
277  },
278  .....
279  },
280  .....
281  },
282  .....
283  },
284  .....
285  },
286  .....
287  },
288  .....
289  },
290  .....
291  },
292  .....
293  },
294  .....
295  },
296  .....
297  },
298  .....
299  },
300  .....
301  },
302  .....
303  },
304  .....
305  },
306  .....
307  },
308  .....
309  },
310  .....
311  },
312  .....
313  },
314  .....
315  },
316  .....
317  },
318  .....
319  },
320  .....
321  },
322  .....
323  },
324  .....
325  },
326  .....
327  },
328  .....
329  },
330  .....
331  },
332  .....
333  },
334  .....
335  },
336  .....
337  },
338  .....
339  },
340  .....
341  },
342  .....
343  },
344  .....
345  },
346  .....
347  },
348  .....
349  },
350  .....
351  },
352  .....
353  },
354  .....
355  },
356  .....
357  },
358  .....
359  },
360  .....
361  },
362  .....
363  },
364  .....
365  },
366  .....
367  },
368  .....
369  },
370  .....
371  },
372  .....
373  },
374  .....
375  },
376  .....
377  },
378  .....
379  },
380  .....
381  },
382  .....
383  },
384  .....
385  },
386  .....
387  },
388  .....
389  },
390  .....
391  },
392  .....
393  },
394  .....
395  },
396  .....
397  },
398  .....
399  },
400  .....
401  },
402  .....
403  },
404  .....
405  },
406  .....
407  },
408  .....
409  },
410  .....
411  },
412  .....
413  },
414  .....
415  },
416  .....
417  },
418  .....
419  },
420  .....
421  },
422  .....
423  },
424  .....
425  },
426  .....
427  },
428  .....
429  },
430  .....
431  },
432  .....
433  },
434  .....
435  },
436  .....
437  },
438  .....
439  },
440  .....
441  },
442  .....
443  },
444  .....
445  },
446  .....
447  },
448  .....
449  },
450  .....
451  },
452  .....
453  },
454  .....
455  },
456  .....
457  },
458  .....
459  },
460  .....
461  },
462  .....
463  },
464  .....
465  },
466  .....
467  },
468  .....
469  },
470  .....
471  },
472  .....
473  },
474  .....
475  },
476  .....
477  },
478  .....
479  },
480  .....
481  },
482  .....
483  },
484  .....
485  },
486  .....
487  },
488  .....
489  },
490  .....
491  },
492  .....
493  },
494  .....
495  },
496  .....
497  },
498  .....
499  },
500  .....
501  },
502  .....
503  },
504  .....
505  },
506  .....
507  },
508  .....
509  },
510  .....
511  },
512  .....
513  },
514  .....
515  },
516  .....
517  },
518  .....
519  },
520  .....
521  },
522  .....
523  },
524  .....
525  },
526  .....
527  },
528  .....
529  },
530  .....
531  },
532  .....
533  },
534  .....
535  },
536  .....
537  },
538  .....
539  },
540  .....
541  },
542  .....
543  },
544  .....
545  },
546  .....
547  },
548  .....
549  },
550  .....
551  },
552  .....
553  },
554  .....
555  },
556  .....
557  },
558  .....
559  },
560  .....
561  },
562  .....
563  },
564  .....
565  },
566  .....
567  },
568  .....
569  },
570  .....
571  },
572  .....
573  },
574  .....
575  },
576  .....
577  },
578  .....
579  },
580  .....
581  },
582  .....
583  },
584  .....
585  },
586  .....
587  },
588  .....
589  },
590  .....
591  },
592  .....
593  },
594  .....
595  },
596  .....
597  },
598  .....
599  },
600  .....
601  },
602  .....
603  },
604  .....
605  },
606  .....
607  },
608  .....
609  },
610  .....
611  },
612  .....
613  },
614  .....
615  },
616  .....
617  },
618  .....
619  },
620  .....
621  },
622  .....
623  },
624  .....
625  },
626  .....
627  },
628  .....
629  },
630  .....
631  },
632  .....
633  },
634  .....
635  },
636  .....
637  },
638  .....
639  },
640  .....
641  },
642  .....
643  },
644  .....
645  },
646  .....
647  },
648  .....
649  },
650  .....
651  },
652  .....
653  },
654  .....
655  },
656  .....
657  },
658  .....
659  },
660  .....
661  },
662  .....
663  },
664  .....
665  },
666  .....
667  },
668  .....
669  },
670  .....
671  },
672  .....
673  },
674  .....
675  },
676  .....
677  },
678  .....
679  },
680  .....
681  },
682  .....
683  },
684  .....
685  },
686  .....
687  },
688  .....
689  },
690  .....
691  },
692  .....
693  },
694  .....
695  },
696  .....
697  },
698  .....
699  },
700  .....
701  },
702  .....
703  },
704  .....
705  },
706  .....
707  },
708  .....
709  },
710  .....
711  },
712  .....
713  },
714  .....
715  },
716  .....
717  },
718  .....
719  },
720  .....
721  },
722  .....
723  },
724  .....
725  },
726  .....
727  },
728  .....
729  },
730  .....
731  },
732  .....
733  },
734  .....
735  },
736  .....
737  },
738  .....
739  },
740  .....
741  },
742  .....
743  },
744  .....
745  },
746  .....
747  },
748  .....
749  },
750  .....
751  },
752  .....
753  },
754  .....
755  },
756  .....
757  },
758  .....
759  },
760  .....
761  },
762  .....
763  },
764  .....
765  },
766  .....
767  },
768  .....
769  },
770  .....
771  },
772  .....
773  },
774  .....
775  },
776  .....
777  },
778  .....
779  },
780  .....
781  },
782  .....
783  },
784  .....
785  },
786  .....
787  },
788  .....
789  },
790  .....
791  },
792  .....
793  },
794  .....
795  },
796  .....
797  },
798  .....
799  },
800  .....
801  },
802  .....
803  },
804  .....
805  },
806  .....
807  },
808  .....
809  },
810  .....
811  },
812  .....
813  },
814  .....
815  },
816  .....
817  },
818  .....
819  },
820  .....
821  },
822  .....
823  },
824  .....
825  },
826  .....
827  },
828  .....
829  },
830  .....
831  },
832  .....
833  },
834  .....
835  },
836  .....
837  },
838  .....
839  },
840  .....
841  },
842  .....
843  },
844  .....
845  },
846  .....
847  },
848  .....
849  },
850  .....
851  },
852  .....
853  },
854  .....
855  },
856  .....
857  },
858  .....
859  },
860  .....
861  },
862  .....
863  },
864  .....
865  },
866  .....
867  },
868  .....
869  },
870  .....
871  },
872  .....
873  },
874  .....
875  },
876  .....
877  },
878  .....
879  },
880  .....
881  },
882  .....
883  },
884  .....
885  },
886  .....
887  },
888  .....
889  },
890  .....
891  },
892  .....
893  },
894  .....
895  },
896  .....
897  },
898  .....
899  },
900  .....
901  },
902  .....
903  },
904  .....
905  },
906  .....
907  },
908  .....
909  },
910  .....
911  },
912  .....
913  },
914  .....
915  },
916  .....
917  },
918  .....
919  },
920  .....
921  },
922  .....
923  },
924  .....
925  },
926  .....
927  },
928  .....
929  },
930  .....
931  },
932  .....
933  },
934  .....
935  },
936  .....
937  },
938  .....
939  },
940  .....
941  },
942  .....
943  },
944  .....
945  },
946  .....
947  },
948  .....
949  },
950  .....
951  },
952  .....
953  },
954  .....
955  },
956  .....
957  },
958  .....
959  },
960  .....
961  },
962  .....
963  },
964  .....
965  },
966  .....
967  },
968  .....
969  },
970  .....
971  },
972  .....
973  },
974  .....
975  },
976  .....
977  },
978  .....
979  },
980  .....
981  },
982  .....
983  },
984  .....
985  },
986  .....
987  },
988  .....
989  },
990  .....
991  },
992  .....
993  },
994  .....
995  },
996  .....
997  },
998  .....
999  },
1000  .....

```

```

'A number indicates the number of cells of the area containing this cell.'],
'block': [
'A chain cannot contain more than one block of one size and shape.',
'Divide all of the board into blocks.',
'Lines must surround a cell and cannot enter a block.',
'The number written on a cell indicates the number of black cells that make up the "block" containing that cell.',
'Blocks are considered to have the same shape even if a block matches the rotated or mirror image of another block in a chain.',
'There can be no blocks with fewer (0 - 4) or more (six and over) cells.',
'A number in a cell shows how many lines of a block there are around it, the number includes the lines of the outer frame.',
'Same sized blocks cannot touch each other, horizontally or vertically.',
'Fill each block with the same number horizontally or vertically.',
'Consecutive black cells are called a "block".',
'Every block must contain one number or one "?".',
'All blocks must be part of a chain.',
'A block must contain a pair of areas of white cells and gray cells having the same form (size and shape).',
'When two or more blocks are in contact at a corner it is called a "chain".',
'A block can contain any number of cells with the number.',
'Divide the grid into blocks by drawing solid lines over the dotted lines.',
'Learn more about this author: Divide the grid into blocks of five cells.',
'Each block contains as many cells as the number in the block.',
'A number indicates the number of cells of that color in the block; the number corresponds to half the number of cells of the block.',
'There may be as many as 5 cells with numbers in a block.',
'If there is a cell with a "?", you do not know the number of black cells that make up the block containing that cell.'],
}

```

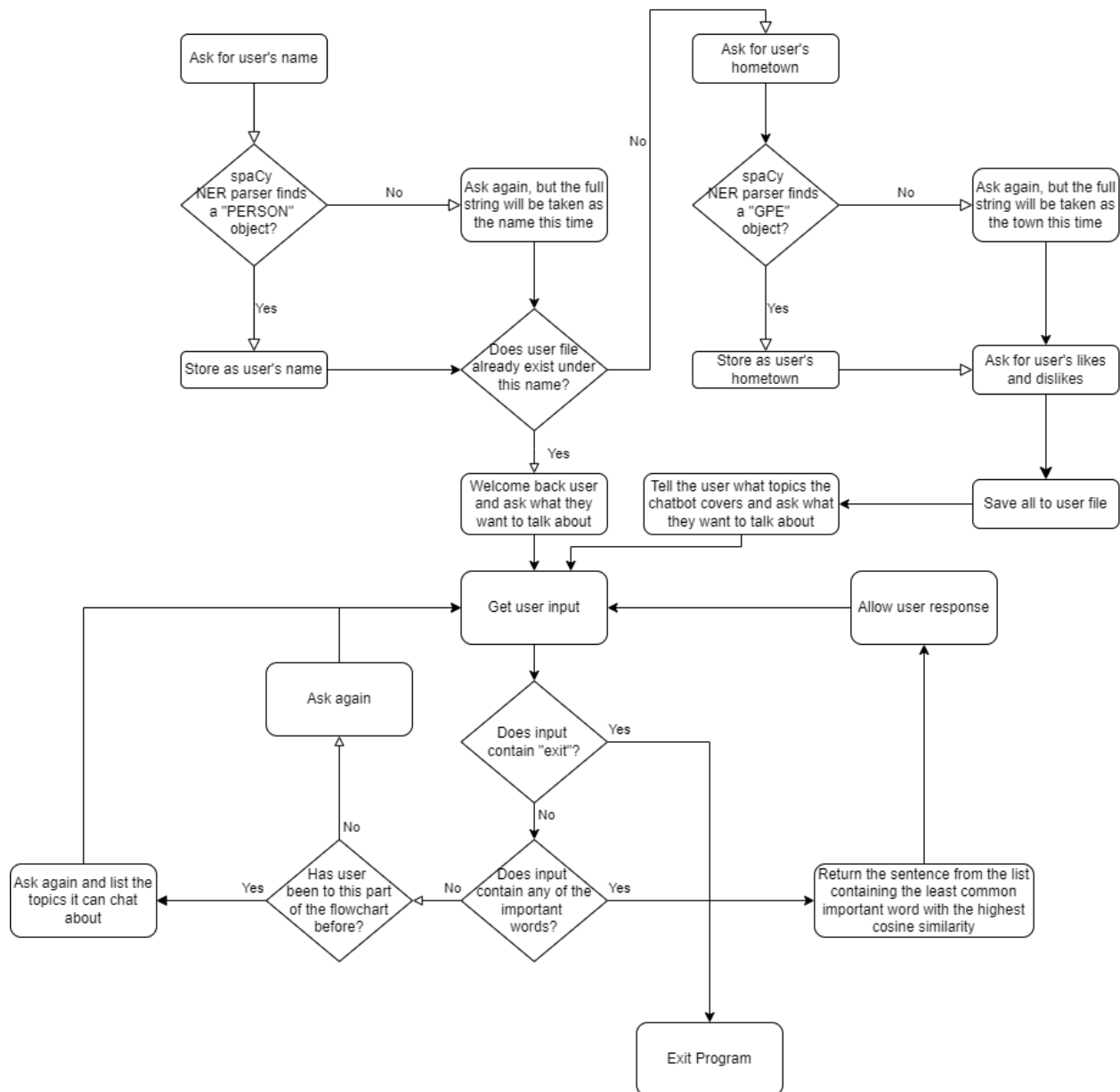
Chatbot

System Description

The chatbot consists of a rules-based beginning, which eventually becomes a query/response lookup with the knowledge base. The rules-based beginning uses named entity recognition (NER) to identify the user's name and their hometown, however it defaults to just using the whole string as some names and locations always get misidentified regardless of phrasing (from testing, I could rarely get it to pick up my own name, "AJ", and there are doubtless others).

Once the query section begins, however, the input from the user is modified using WordNetLemmatizer and NLTK, tokenizing, lowercasing and lemmatizing the words, as well as removing stopwords. The input is then checked against the 10 critical terms—if the input contains none of them, the chatbot will identify the topics it knows and ask the user to repeat. The critical words are listed in order of decreasing frequency and the least common word is used as a base. This operates off the assumption that if a sentence contains a less common word, it is probably based around that sentence. So the chatbot looks at all the sentences from the knowledge base and compares them using cosine similarity to the query sentence. The most similar sentence is then returned to the user, and the user is allowed to ask another question.

Flowchart



Sample Dialogue:

First time user:

Hello, what is your name?
 >Jonathan
 [Jonathan]
 Nice to meet you! What is your hometown?
 >Phoenix
 [Phoenix]

Alright! What is something you like?
>Video Games
Okay, what is something you dislike?
>Twitter
Nice to meet you, Jonathan! I am a chatbot with information on the Japanese puzzle publisher, Nikoli, which is most well known for creating Sudoku.
What would you like to talk about? (type "exit" to quit.)
>Have they made video games?
Sorry, I don't understand. I can talk about the company Nikoli, their puzzles (including Sudoku), their rules and how to solve them, Nikoli's books and magazines, or common puzzle elements like cells and blocks.
>What is the best game of sudoku ever played?
"What's next after the sudoku craze?" People are looking to Nikoli for answers.
>exit

Returning user:

Hello, what is your name?
>AJ
Sorry, I didn't quite get that. What is your name?
>AJ
Oh, hello, AJ! I remember you--you're the person from Phoenix who likes Food and dislikes Bugs.
What would you like to talk about? (type "exit" to quit.)
>Where can I buy Nikoli's magazines?
"Puzzle Communication Nikoli Vol.186" Nikoli is the most advanced pencil puzzle magazine in Japan.
>where can I buy nikoli's books?
"Puzzle Communication Nikoli Vol.186", "Nikoli Puzzle Book Reprint Edition Kakuro 2" and "Nikoli Puzzle Book Reprint Edition Slitherlink 2" are now on sale.
>how might a puzzle use cells?
Balloons are light and float, so they must be placed in one of the cells at the top, or in a cell right under a black cell (filled-in cell) or under other balloons.
>exit

As you can see, the parser could not pick up "AJ" as a name.

Hello, what is your name?
>hi my name is matthew
Sorry, I didn't quite get that. What is your name?
>My name is Matthew
Nice to meet you! What is your hometown?
>Murietta California
Alright! What is something you like?
>I like to play Sudoku on my phone.
Okay, what is something you dislike?

>Cold food

Nice to meet you, My name is Matthew! I am a chatbot with information on the Japanese puzzle publisher, Nikoli, which is most well known for creating Sudoku.

What would you like to talk about? (type "exit" to quit.)

>I want to learn about Nikoli

With Nikoli logic puzzles you need just paper, pencil, and an active mind, no expensive programs or training..... nikoli Puzzles page nikoli book shop "Puzzle Communication Nikoli Vol.186", "Nikoli Puzzle Book Reprint Edition Kakuro 2" and "Nikoli Puzzle Book Reprint Edition Slitherlink 2" are now on sale.

>Tell me more about the creation of sudoku

Each puzzle is handcrafted (not computer generated, like most generic sudoku books) by the experts at Nikoli Publishing, the Japanese creators of the game, so solving one is like pitting your mind against that of a sudoku master.

>Tell me about other games published by Nikoli

With Nikoli logic puzzles you need just paper, pencil, and an active mind, no expensive programs or training..... nikoli Puzzles page nikoli book shop "Puzzle Communication Nikoli Vol.186", "Nikoli Puzzle Book Reprint Edition Kakuro 2" and "Nikoli Puzzle Book Reprint Edition Slitherlink 2" are now on sale.

>Tell me about the founder of the company

Sorry, I don't understand. I can talk about the company Nikoli, their puzzles (including Sudoku), their rules and how to solve them, Nikoli's books and magazines, or common puzzle elements like cells and blocks.

>Tell me about the founder of Nikoli

With Nikoli logic puzzles you need just paper, pencil, and an active mind, no expensive programs or training..... nikoli Puzzles page nikoli book shop "Puzzle Communication Nikoli Vol.186", "Nikoli Puzzle Book Reprint Edition Kakuro 2" and "Nikoli Puzzle Book Reprint Edition Slitherlink 2" are now on sale.

>exit

AJs-MacBook-Air:Spring 2024 NLP AJ\$ /usr/local/bin/python3 "/Users/AJ/VC/Spring 2024 NLP/Chatbot/chatbot.py"

Hello, what is your name?

>My name is Matthew

Nice to meet you! What is your hometown?

>^CKeyboardInterrupt

AJs-MacBook-Air:Spring 2024 NLP AJ\$ /usr/local/bin/python3 "/Users/AJ/VC/Spring 2024 NLP/Chatbot/chatbot.py"

Hello, what is your name?

>matthew

Sorry, I didn't quite get that. What is your name?

>My name is Matthew

Oh, hello, My name is Matthew! I remember you--you're the person from California who likes I like to play Sudoku on my phone. and dislikes Cold food.

What would you like to talk about? (type "exit" to quit.)

>

The NER has trouble recognizing lowercase names, and is not set up to handle parsing likes/dislikes.

Knowledge base:

The format of the knowledge base is a dictionary, with the important word as the key and a list of sentence options as the value. The entire knowledge base is over 200 lines long, so I will not paste it here, but the format is as follows:

```
{
  'nikoli': [
    'nikoli book shop page "Puzzle Communication Nikoli Vol.186" is now on sale.',
    'All our magazines and books are available at the "nikoli book shop" page.',
    ...
  ],
  'puzzle': [
    'The joy of solving puzzles is more than a passing fad, it enriches lives.',
    ...
  ],
  ...
  {Sections for all important words}
  ...
  'block': [...]
```

User profiles:

The format of the user profiles is a simple dictionary in the following format:

AJ.p

```
{ "hometown": "Phoenix",
  "likes": "Food",
  "dislikes": "Bugs" }
```

This format could be easily extended to include more information about the user, however this is all this iteration of the program stores.

Analysis

Since this chatbot was limited to being rules-based as well as some other challenging restraints (must utilize a relatively small knowledge base scraped from the internet, this knowledge base must largely be cleaned up by the computer rather than by a human, and created within a relatively short development cycle by one person), the

output is not very good or convincing. Some inquiries it succeeds with (it seems to be the most successful on asking where to play puzzles or where to buy books/magazines), whereas others it often outputs nonsense or incorrect information (such as asking for the rules of a specific puzzle). A more robust system of detecting likes and dislikes using sentiment analysis would also be ideal—however developing a system to differentiate between sentences that state a preference, a query, or a combination of both was a bit out of scope for the scale of this project. More development time and a more carefully curated/manually written knowledge base likely could improve the quality of this system, however this would no longer fit the mechanical requirements of this specific project.

One particular challenge was picking suitable important words—it is easy to tell when a word is important, but it's another challenge to try and pick words that will output useful information, and I tried to pick the best out of the 50 most relevant TF-IDF words. However, some words are more prone to outputting nonsense or nonsequiturs—the “block”, “cells”, and “Japan” words being the worst offenders, with the first two missing most of their context most of the time, and the second appearing incidentally a lot and returning information that may not be all that relevant.

For outside evaluation, I asked my friends Jonathan, Ethan, and Matthew and charted their responses in this table.

	Strongly disagree	Partially disagree	Neither agree nor disagree	Partially agree	Strongly agree
The chatbot answered my questions accurately.					
		Jonathan, Ethan	AJ, Matthew		
The chatbot remembered who I was on a second run.					
		Matthew			AJ, Jonathan, Ethan
The chatbot gave reasonable responses to my queries.					

		AJ	Matthew	Jonathan, Ethan	
The chatbot was a convincing conversation partner.					
	AJ, Jonathan, Matthew		Ethan		

Overall, the replies are only somewhat relevant to the queried topic. A larger database, a more robust data cleaning process, and word similarity checking rather than just simple word-by-word cosine similarity might help diminish these issues. The strongest part is the more rules-based beginning, but that too occasionally runs into issues if the user types in all lowercase or if they phrase their likes and dislikes as sentences.