

```

1 from __future__ import division
2 from __future__ import print_function
3 import random
4 import gensim
5 from random import randint
6 from translate import Translator
7 import numpy as np
8 from transliterate import translit
9 import argparse, sys
10 import matplotlib.pyplot as plt
11 from scipy.spatial import distance
12 import seaborn as sns
13 import pandas as pd
14 from nltk.corpus import wordnet
15 from tabulate import tabulate
16 def checkStoredWords(kwords, word):
17     """
18     This function updates a list of known words with a new word.
19     If the spell type and language exists in the list the value is append
20     by 1 otherwise, it is appended to the end of the list with a value of
21     1.
22     :param kwords: List of spell types and language with associated
23     frequencies.
24     :param word: One being the spell type and the other being the
25     origin language.
26     :type kwords: [[str, str], int]...
27     :type word: str
28     :return: the updated list of known words.
29     """
30     found = False
31     for kword in kwords:
32         if kword[0] == word:
33             kword[1] += 1
34             found = True
35     if found == False:
36         kwords.append([word, int(1)])
37     return kwords
38
39 def count_instances(fname):
40     """
41     Reads supplied file, where it splits it up. Then it appends
42     each word to the data set building a list of words and frequencies
43     using checkStoredWords(kwords, word).
44     :param fname: This is the name of the CSV file in which the spell
45     data is stored.
46     :type fname: str
47     :return: returns a list of languages and the probabilities for
48     each one.
49     """

```

```

48     file = open(fname, 'r')
49     data = []
50
51     for line in file:
52         temp = line.rstrip()
53         temp = temp.split(",")
54         data = checkStoredWords(data, temp)
55     file.close()
56     data = calcProb(data)
57     return data
58
59 def totalSpells(data):
60     """
61         Counts the number of spells in the dataset.
62
63         :param data: List of spell types and origin language with
64         frequency.
65         :type data: [[[str,str], int]...]
66         :return: an integer value of total number of spells.
67     """
68     total = 0
69     for d in data:
70         total += d[1]
71     return total
72
73 def calcProb(data):
74     """
75         Calculates the probabilities for spells of each type.
76
77         :param data: List of spell types and origin language with
78         frequency.
79         :type data: [[[str, str], int]...]
80         :return: A list of type of spells and their associated
81         probabilities.
82     """
83     total = totalSpells(data)
84     prob = 0.0
85     for d in data:
86         prob = d[1] / total
87         d.append(prob)
88     return data
89
90 def generateScale(data):
91     """
92         This stacks the probabilities of spells so that each spell
93         has a boundary in which it a spell can be selected over another.
94
95         :param data: list of spell names and their associated frequencies
96         and probabilities.

```

```

95         :type data: [[[str,str],int,float]...]
96         :return: a list of spells and the value between 0-1 in which that
97         name will be selected.
98         """
99         value = 0
100        index = -1
101        scale = []
102        for d in data:
103            value += d[2]
104            index += 1
105            scale.append((value, d[0]))
106        return scale
107
108    def getSpellType(scale, rndNum):
109        """
110            Selects a spell according to the random number passed.
111
112            :param scale: A list of tuples which contains the probability
113            associated with each spell and type.
114            :param rndNum: The random number used to select a spell type.
115            :type scale: [(str,str,float)..]
116            :type rndNum: float
117            :return: A string which is the spell type.
118            """
119        for i in range(-1, len(scale) - 1):
120            if i == -1:
121                temp2 = scale[i + 1]
122                if rndNum >= 0:
123                    if rndNum < temp2[0]:
124                        return temp2[1]
125            else:
126                temp = scale[i]
127                temp2 = scale[i + 1]
128                if rndNum >= temp[0]:
129                    if rndNum < temp2[0]:
130                        return temp2[1]
131
132        temp2 = scale[0]
133        return temp2[1]
134
135    def is_valid(string):
136        """
137            check to see whether a word consists of alpha characters.
138
139            :param string: The string to be checked.
140            :type string: str
141            :return: Boolean value.
142            """
143        if string.isalpha():
144            return False
145        return True

```

```

142
143
144 def langCode(language): #this now works with python 2.7 i believe.
145     """
146     Converts a language name into a language code for the translator.
147
148     :param language: Full name of the language, for example latin.
149     :type language: tr
150     :return: The string code for the language.
151     """
152     return {
153         'Latin': 'la',
154         'Greek': 'el',
155         'Portuguese': 'pt',
156         'West African Sidiki': 'it', # CANT BE TRANSLATED. - Returns
157         italian
158         'Aramaic': 'el', # CANT BE TRANSLATED - RETURNS GREEK
159         'Pig Latin': 'PL', # implement a seperate function to
160         convert to pig latin.
161         'English': 'en',
162         'French': 'fr',
163         'Spanish': 'es',
164         'Italian': 'it',
165     }.get(language, 'la') # returns latin as default - if language
166     is not found.
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188

```

```

189     logfile.write(text.encode("utf-8") + "\n")
190     logfile.close()
191
192     def sentenceToWord(sentence, model, oword):
193         """
194         Takes a string and converts it into a vector. Then from that it
195         picks a similar word that doesn't contain an underscore.
196
197         :param sentence: A string which contains a sentence to be
198         converted into one word.
199         :type sentence: str
200         :return: A string containing a similar word.
201         """
202
203         sentence = sentence.split()
204         output = []
205         top_val = 20
206         selected = []
207         bogus_words = 0
208         for word in sentence:
209             try:
210                 output.append(model[word])
211             except KeyError:
212                 log("key error in vector file" + word)
213
214         output = np.array(output)
215         vector_sum = output.sum(axis=0)
216         output = model.most_similar(positive=[vector_sum], topn=top_val)
217         final_output = output[randint(0, (top_val - 1))]
218         while is_valid(final_output[0]):
219             num = randint(0, top_val - 1)
220             final_output = output[num]
221             if num in selected:
222                 if len(selected) == top_val:
223                     top_val = top_val * 2
224                     output = model.most_similar(positive=[vector_sum],
225                                                 topn=top_val)
226             else:
227                 selected.append(num)
228
229         bogus_words+=1
230         return final_output, bogus_words
231
232     def pigLatin(source):
233         """
234         Takes a source string and converts it from english to pig
235         latin.
236
237         :param source: Takes string of english words and changes it into
238         pig latin.
239         :type source: str
240         :return: a string containing pig latin words.
241         """

```

```

236 letters = ['sh', 'gl', 'ch', 'ph', 'tr', 'br', 'fr', 'bl', 'gr',
237 'st', 'sl', 'cl', 'pl', 'fl']
238 source = source.split()
239 for k in range(len(source)):
240     i = source[k]
241     if i[0] in ['a', 'e', 'i', 'o', 'u']:
242         source[k] = i + 'ay'
243     elif f(i) in letters:
244         source[k] = i[2:] + i[:2] + 'ay'
245     elif i.isalpha() == False:
246         source[k] = i
247     else:
248         source[k] = i[1:] + i[0] + 'ay'
249 return ' '.join(source)
250
251 def f(str):
252     """
253     Returns the first two chacters from the string.
254
255     :param str: A word that is passed.
256     :type str: str
257     :return: a string that only contains the first two letters.
258
259     """
260     if len(str) == 1:
261         return str[0]
262     return str[0] + str[1]
263
264 def generateSpell(sentence, model, oword):
265     """
266     Generates a Spell from a sentence.
267
268     :param sentence: string which is the definition of the spell you
269     want to create.
270     :type sentence: str
271     :return: list containing the spell and the spell type.
272     :param model: loaded vector orepresentation of words.
273     :type model: data file loaded.
274
275     """
276     spell = []
277     vector, temp_bogus = sentenceToWord(sentence, model, oword)
278     vector = vector[0]
279     scale = generateScale(count_instances('spell_prob.csv'))
280     selection = random.random()
281     spell_meta = getSpellType(scale, selection)
282
283     try:
284         target_lang = langCode(spell_meta[1])
285     except:
286         log("langCode function didn't work. Using default latin.")
287         target_lang = "la"
288
289     if target_lang == "PL":
290         spell.append(pigLatin(vector))
291     else:

```

```

283         spell.append(translate2(vector, target_lang))
284         spell.append(spell_meta[0])
285         spell.append(vector) #The original word before translation is
also added.
286         return spell, temp_bogus
287
288
289 def load_vectors(path, is_binary):
290     """
291     This loads the vectors supplied by the path.
292
293     :param path: The path to the vector file
294     :type path: str
295     :param is_binary: states whether file is a binary file.
296     :type is_binary: boolean
297     :return: The loaded model.
298     """
299     print("Loading: ", path)
300     model = gensim.models.Word2Vec.load_word2vec_format(path,
301                                                         binary=is_binary)
302     model.init_sims(replace=True)
303     print("Loaded: ", path)
304     return model
305
306
307 def is_synonym(n_word, o_word):
308     """
309     This function uses a combination of NLTK's wordnet to
list all synonyms for a word and to check if a new word is a
synonym.
310
311     :param n_word: The new word generated.
312     :type n_word: str
313     :param o_word: The original word in the definition.
314     :type o_word: str
315     :return: Returns a boolean indicating whether n_word is a synonym
316     of o_word.
317     """
318     synonyms=[]
319     synsets = wordnet.synsets(o_word)
320     for synset in synsets:
321         synonyms = synonyms+ synset.lemma_names()
322
323     return n_word in synonyms
324
325
326 def run_experiment(model, num_experiments):
327     """
328     This function runs the experiments with the paramters set.
It then returns all the necessary data for processing and
output.
329
330     :param model: The vectors loaded.
331     :type model: The loaded vector object
332     :param num_experiments: The number of experiments to run.

```

```

330         :type num_experiments: int
331         :return: A list of averages scores, one entry per experiment.
332         :return: A list of the average number of synonyms produced,
one entry per experiment.
333         :return: The average score across the experiments.
334         :return: A list of average cosine similarity scores, one
entry per experiment.
335         :return: The number of experiments.
336         :return: A list containing the number of bogus words
produced, one entry per experiment.
337         :return: A list containing lists with each sublist containing
the scores produced for that definition length.
338         :return: A list containing list with each sublist containing
number of bogus words produced for that definition length.
341         """
342         average = 0.0
343         iterationCount = 0
344         scores = []
345         cos_dists = []
346         avg_cos_dists = []
347         syn_experiments = []
348         bword_counts = []
349         scores_per_spell=[[] for x in range(10)] #size of definition
length.
350         table1 = []
351         table2 = []
352         bwords_spell= [[] for x in range(10)] #size of definition length.
353         for i in range(0, num_experiments):
354             table1 = []
355             table2 = []
356             print("-----", i, "-----")
357             log("-----"+str(i) + "-----")
358             bogus_words = 0
359             spellFile = open("spells.csv")
360             entry = []
361             score = 0
362             count = 0
363             syn_counts = 0
364             for line in spellFile:
365                 count+=1
366                 line = line.strip("\n")
367                 entry = line.split(",")
368                 spell, temp_bogus = generateSpell(entry[1],
model,entry[3])
369                 bwords_spell[len(entry[1].split(" "))].append(temp_bogus)
370                 bogus_words+= temp_bogus
371
372                 if args.verbose:
373                     print("Your new spell is: ", spell[0])
374
375                 if spell[2].lower() not in entry[1].split():
376                     score +=1
377                     scores_per_spell[len(entry[1].split(" "))].append(1)
378                 else:
379                     scores_per_spell[len(entry[1].split(" "))].append(0)
380
381             table1.append([spell[0]])
382             table2.append([spell[2]])

```



```

377         #calculate the cosine similarity.
378         og_wd = model[entry[-1].strip()]
379         nw_wd = model[spell[-1]]
380         cos_dists.append(distance.cosine(og_wd, nw_wd))
381
382         if is_synonym(spell[2].lower(), entry[-1]):
383             syn_counts +=1
384
385     print("Experiment Results")
386     print("Num of spells that don't feature in definition: ",
387 score)
388     print("Percentage: ", ((float(score)/count) * 100), "%")
389     print("Average Cosine-similarity:", float(sum(cos_dists) /
len(cos_dists)))
390     print("Num of spells which are synonyms: ", syn_counts)
391     print("Num of words selected that are not real words: ",
bogus_words)
392     scores.append((float(score)/count) * 100)
393     syn_experiments.append(syn_counts)
394     bword_counts.append(bogus_words)
395     spellFile.close()
396     iterationCount +=1
397     average += (float(score)/count)*100
398     avg_cos_dists.append(float(sum(cos_dists) / len(cos_dists)))
399     return scores, syn_experiments, average, avg_cos_dists,
iterationCount, bword_counts, scores_per_spell, bwords_spell
400
401 #
402 # Main part of the program.
403 #
404 #
405 if __name__ == '__main__':
406     parser = argparse.ArgumentParser(
407         'Use Word2Vec or GloVe datasets to generate Harry Potter
Spells')
408     parser.add_argument('--glove', action='store_const',
409         const = 'glove',
410         help='Use the GloVe dataset instead of the default
Word2Vec.')
411     parser.add_argument('--exp',
412         help="Specifies the number of experiments on this run.
Default is 20.",
413         action='store', type=int)
414     parser.add_argument('--verbose', action='store_const',
415         const = 'verbose',
416         help='Prints out the spell names')
417     parser.add_argument('--comp', action= 'store_const',
418         const='comp',
419         help = "Runs the word2vec vectors, and the GloVe
vectors")
420     args = parser.parse_args()
421
422     logFile = open("log.txt", 'w' )
423     logFile.close()

```

```

424     num_experiments = 20
425
426     if args.exp != None:
427         num_experiments = args.exp
428
429     if args.comp: # comparison mode.
430         print("Compare Mode")
431         log("-----Compare
Mode-----")
432         print("Vectors used: Word2Vec")
433         log("-----+ Vectors used: Word2Vec"+
434             "-----")
435         model = load_vectors(
436             True)
437             ".../..vectors/GoogleNews-vectors-negative300.bin",
438
439         #Run word2vec experiments and then stores data in dataframe.
440         w_scores, w_syn_experiments, w_average, w_avg_cos_dists,
441         iterationCount, w_bword_counts, w_spells_per, w_bwords_per =
442         run_experiment(
443             model, num_experiments)
444         w_vec=["word2vec" for x in w_scores]
445         del model
446         print("Vectors used: GloVe")
447         log("-----+ Vectors used: GloVe"+
448             "-----")
449         model = load_vectors(".../..vectors/glove.txt.vw", False)
450
451         # run experiments and move results into data frame.
452         g_scores, g_syn_experiments, g_average, g_avg_cos_dists,
453         iterationCount, g_bword_counts, g_spells_per, g_bwords_per=
454         run_experiment(model, num_experiments)
455         g_vec = ["glove" for x in g_scores]
456
457         scores=w_scores + g_scores
458         syn_experiments = w_syn_experiments + g_syn_experiments
459         avg_cos_dists = w_avg_cos_dists + g_avg_cos_dists
460         bword_counts = w_bword_counts + g_bword_counts
461         vectors = w_vec + g_vec
462
463         ##for the ts plots
464         g_vec = ["GloVe" for x in g_spells_per]
465         w_vec = ["Word2Vec" for x in w_spells_per]
466         bwords_per = w_bwords_per + g_bwords_per
467         spells_per = w_spells_per + g_spells_per
468         vec = w_vec + g_vec
469
470         ##adds values for empty rows.
471         for row in spells_per:
472             if len(row) == 0:
473                 row.append(0)
474         for row in bwords_per:
475             if len(row) == 0:
476                 row.append(0)
477
478         spells_per_avg = [float(sum(l)/len(l)) for l in spells_per]
479         length= [x for x in range(1, len(w_spells_per)+1)] + [x for x
in range(

```

```

471         1, len(g_spells_per)+1]]
472     bwords_per_avg = [float(sum(l)/len(l)) for l in bwords_per]
473     len_results = pd.DataFrame({"originality":spells_per_avg,
474                                "length":length, "bwords":bwords_per_avg, "vectors":vec})
475
476     box_len = []
477     box_score= []
478     box_vec=[]
479
480     for i in range(0, len(w_spells_per)):
481         for row in w_spells_per[i]:
482             box_len.append(i+1)
483             box_score.append(row)
484             box_vec.append("word2vec")
485
486         for row2 in g_spells_per[i]:
487             box_len.append(i+1)
488             box_score.append(row2)
489             box_vec.append("GloVe")
490
491     box_data = pd.DataFrame({"length":box_len,
492                              "originality":box_score, "vectors":box_vec})
493
494     #originality vs size plots.
495     ax = sns.tsplot(time="length", value="originality",
496                     unit="vectors", condition="vectors",
497                     data=len_results )
498     sns.plt.show()
499     ax = sns.distplot(box_score)
500     sns.plt.show()
501
502     #box plot
503     ax = sns.boxplot(x="length", y = "originality",
504                      hue="vectors",
505                      data=box_data)
506     sns.plt.show()
507
508     box_len = []
509     box_score= []
510     box_vec=[]
511
512     for i in range(0, len(w_bwords_per)):
513         for row in w_bwords_per[i]:
514             box_len.append(i+1)
515             box_score.append(row)
516             box_vec.append("word2vec")
517
518         for row2 in g_bwords_per[i]:
519             box_len.append(i+1)
520             box_score.append(row2)
521             box_vec.append("GloVe")
522
523     box_data = pd.DataFrame({"length":box_len,
524                              "bwords":box_score,
525                              "vectors":box_vec})
526
527     #histogram
528     ax = sns.distplot(box_score)

```

```

518     sns.plt.show()
519     #box plot
520     ax = sns.boxplot(x="length", y = "bwords", hue="vectors",
521                     data=box_data)
522     sns.plt.show()
523     ##output results.
524     print("-----word2vec Experiment
Results-----")
525     print("The mean average percentage over ", iterationCount ,
526           "tests: ", (w_average/iterationCount), "%")
527     print("The mean cosine simalarity over ", iterationCount,
528           "tests: ",
529           float(sum(w_avg_cos_dists)/ len(w_avg_cos_dists)))
530     print("The mean amount of synonyms",
531           (sum(w_syn_experiments)/ iterationCount))
532     print("Average number of words that are not fit for
translation: ",float(sum(w_bword_counts)/iterationCount))
533
534     print("-----GloVe Experiment
Results-----")
535     print("The mean average percentage over ", iterationCount ,
536           "tests: ",(g_average/iterationCount), "%")
537     print("The mean cosine simalarity over ", iterationCount,
538           "tests: ", float(sum(g_avg_cos_dists)/
539 len(g_avg_cos_dists)))
540     print("The mean amount of synonyms",
541           (sum(g_syn_experiments)/ iterationCount))
542     print("Average number of words that are not fit for
translation: ", float(sum(g_bword_counts)/iterationCount))
543
544     results = pd.DataFrame({"scores":scores,
545                             "similarity":avg_cos_dists,
546                             "synonyms":syn_experiments, "vectors":vectors,
547                             "bwords":bword_counts})
548
549     sim = sns.violinplot(x="vectors", y="similarity",
550                          data=results)
551     sns.plt.title("Comparison of Similarity over "+str(
552 iterationCount)+ " experiments")
553     sns.plt.show()
554     sc = sns.violinplot(x="vectors", y="scores", data=results)
555     sns.plt.title("Comparison of accuracy scores over "+str(
556 iterationCount)+ " experiments")
557     sns.plt.show()
558     bw = sns.violinplot(x="vectors", y="bwords", data=results)
559     sns.plt.title("Comparison of invalid words over "+ str(
560 iterationCount)+ " experiments")
561     sns.plt.show()
562     syn = sns.violinplot(x="vectors", y="synonyms", data=results)
563     sns.plt.show()
564     else: # test an individual mode.

```

```

565         if args.glove:
566             print("Vectors used: GloVe")
567             log("-----" + "Vectors used: GloVe"+
568 "-----")
569             model = load_vectors(".././vectors/glove.txt.vw", False)
570         else:
571             print("Vectors used: Word2Vec")
572             log("-----"+ "Vectors used: Word2Vec"+
573 "-----")
574             model = load_vectors(
575                 ".././vectors/GoogleNews-vectors-
576 negative300.bin", True)
577
578         scores, syn_experiments, average, avg_cos_dists,
579 iterationCount, bword_counts, spells_per, bwords_per=
580 run_experiment(model, num_experiments)
581         print("-----Experiment Results-----")
582         print("The mean average percentage over ", iterationCount ,
583             "tests: ", (average/iterationCount), "%")
584         print("The mean cosine similarity over ", iterationCount,
585             "tests: ", float(sum(avg_cos_dists)/
586 len(avg_cos_dists)))
587         print("The mean amount of synonyms",
588             (sum(syn_experiments)/ iterationCount))
589         print("Average number of words that are not fit for
590 translation: ",
591             float(sum(bword_counts)/iterationCount))
592         results = pd.DataFrame({'scores': scores,
593             'similarity': avg_cos_dists})
594
595         #loop through and add an entry to any empty fields.
596         for row in spells_per:
597             if len(row) == 0:
598                 row.append(0)
599
600         spells_per_avg = [float(sum(l)/len(l)) for l in spells_per]
601         length= [x for x in range(0, len(spells_per_avg))]
602
603         vec = ["vector" for x in spells_per_avg]
604         sns.plt.show()
605         ax2 = sns.violinplot(x=results["similarity"])
606         sns.plt.show()
607         ax = sns.violinplot(x="scores", y="similarity", data=results)
608         sns.plt.show()
609
610
611

```