```python
from __future__ import division
from __future__ import print_function
import random
import gensim
from random import randint
from translate import Translator
import numpy as np
from transliterate import translit
import argparse, sys
import matplotlib.pyplot  as plt
from scipy.spatial import distance
import seaborn as sns
import pandas as pd
from nltk.corpus import wordnet
from tabulate import tabulate
def checkStoredWords(kwords, word):
    """

        This function updates a list of known words with a new word. If the spell type and
    language exists in the list the value is append by 1 otherwise, it is appended to
    the end of the list with a value of 1.

    :param kwords: List of spell types and language with associated frequencies.
    :param word: One being the spell type and the other being the origin language.
    :type kwords: [[[str, str], int]...]
    :type word: str
    :return: the updated list of known words.
    """

    found = False
    for kword in kwords:
        if kword[0] == word:
            kword[1] += 1
            found = True
    if found == False:
        kwords.append([word, int(1)])
    return kwords



def count_instances(fname):
    """

        Reads supplied file, where it splits it up. Then it appends each word to the data
    set building a list of words and frequencies using checkStoredWords(kwords, word).

    :param fname: This is the name of the CSV file in which the spell data is stored.
    :type fname: str
    :return: returns a list of languages and the probabilities for each one.
    """

    file = open(fname, 'r')
    data = []

    for line in file:
        temp = line.rstrip()
        temp = temp.split(",")
        data = checkStoredWords(data, temp)
    file.close()
    data = calcProb(data)
    return data


def totalSpells(data):
    """

        Counts the number of spells in the dataset.

    :param data: List of spell types and origin language with frequency.
    :type data: [[[str,str], int]...]
    :return: an integer value of total number of spells.
    """

    total = 0
    for d in data:
        total += d[1]
    return total


def calcProb(data):
    """

        Calculates the probabilities for spells of each type.
```

```python
70          :param data: List of spell types and origin language with frequency.
71          :type data: [[[str, str], int]...]
72          :return: A list of type of spells and their associated probabilities.
            """

74          total = totalSpells(data)
75          prob = 0.0
76          for d in data:
77                  prob = d[1] / total
                    d.append(prob)
78          return data


80  def generateScale(data):
81          """

83              This stacks the probabilities of spells so that each spell has a boundary in which
84          it a spell can be selected over another.

85          :param data: list of spell names and their associated frequencies and probabilities.
86          :type data: [[[str,str],int,float]...]
87          :return: a list of spells and the value between 0-1 in which that name will be selected.
            """

89          value = 0
90          index = -1
91          scale = []
92          for d in data:
93              value += d[2]
94              index += 1
                scale.append((value, d[0]))
95          return scale


97  def getSpellType(scale, rndNum):
98          """

100             Selects a spell according to the random number passed.

101         :param scale: A list of tuples which contains the probability associated with each spell and type.
102         :param rndNum: The random number used to select a spell type.
103         :type scale: [(str,str,float)..]
104         :type rndNum: float
105         :return: A string which is the spell type.
            """

107         for i in range(-1, len(scale) - 1):
108             if i == -1:
109                 temp2 = scale[i + 1]
110                 if rndNum >= 0:
111                     if rndNum < temp2[0]:
                            return temp2[1]
112             else:
113                 temp = scale[i]
114                 temp2 = scale[i + 1]
                    if rndNum >= temp[0]:
115                     if rndNum < temp2[0]:
116                         return temp2[1]

117         temp2 = scale[0]
118         return temp2[1]


120 def is_valid(string):
121         """
122     check to see whether a word consists of alpha characters.

124     :param string: The string to be checked.
125     :type string: str
126     :return: Boolean value.
            """
127         if string.isalpha():
128             return False
            return True




132 def langCode(language): #this now works with python 2.7 i believe.
133         """

135     Converts a language name into a language code for the translator.

136     :param language: Full name of the language, for example latin.
137     :type language: tr
138     :return: The string code for the language.
```

```python
        """
        return {
            'Latin': 'la',
            'Greek': 'el',
            'Portuguese': 'pt',
            'West African Sidiki': 'it',  # CANT BE TRANSLATED. - Returns italian
            'Aramaic': 'el',  # CANT BE TRANSLATED - RETURNS GREEK
            'Pig Latin': 'PL',  # implement a seperate function to convert to pig latin.
            'English': 'en',
            'French': 'fr',
            'Spanish': 'es',
            'Italian': 'it',
        }.get(language, 'la')  # returns latin as default - if language is not found.


def translate2(word, lang):
        """
            Translates a word to a target language.

        :param word: The word you want to convert.
        :param lang: the lang code of the language you want to convert to.
        :type word: str
        :type lang: str
        :return: a string containing the translated word in the latin alphabet.
        """
        translator = Translator(to_lang=lang)
        try:
            out = translator.translate(word)
            if lang == 'el':
                return translit(word, lang, reversed=True)
            return out
        except:
            log("Error Cannot translate: " + word)

def log(text):
        logfile = open("log.txt", "a")
        logfile.write(text.encode("utf-8") + "\n")
        logfile.close()

def sentenceToWord(sentence, model, oword):
        """

        Takes a string and converts it into a vector. Then from that it picks a similar word that doesn't
        contain an underscore.

        :param sentence: A string which contains a sentence to be converted into one word.
        :type sentence: str
        :return: A string containing a similar word.
        """

        sentence = sentence.split()
        output = []
        top_val = 20
        selected = []
        bogus_words = 0
        for word in sentence:
            try:
                output.append(model[word])
            except KeyError:
                log("key error in vector file" + word)

        output = np.array(output)
        vector_sum = output.sum(axis=0)
        output = model.most_similar(positive=[vector_sum], topn=top_val)
        final_output = output[randint(0, (top_val - 1))]
        while is_valid(final_output[0]):
            num = randint(0, top_val - 1)
            final_output = output[num]
            if num in selected:
                if len(selected) == top_val:
                    top_val = top_val * 2
                    output = model.most_similar(positive=[vector_sum], topn=top_val)
            else:
                selected.append(num)

            bogus_words+=1
    # print(final_output[0])
        return final_output, bogus_words


def pigLatin(source):
        """

            Takes a source string and converts it from english to pig latin.

        :param source: Takes string of english words and changes it into pig latin.
```

```python
        :type source: str
        :return: a string containing pig latin words.

        """

        letters = ['sh', 'gl', 'ch', 'ph', 'tr', 'br', 'fr', 'bl', 'gr', 'st', 'sl', 'cl', 'pl', 'fl']
        source = source.split()
        for k in range(len(source)):
            i = source[k]
            if i[0] in ['a', 'e', 'i', 'o', 'u']:
                source[k] = i + 'ay'
            elif f(i) in letters:
                source[k] = i[2:] + i[:2] + 'ay'
            elif i.isalpha() == False:
                source[k] = i
            else:
                source[k] = i[1:] + i[0] + 'ay'
        return ' '.join(source)


def f(str):
    """
    Returns the first two chacters from the string.

    :param str: A word that is passed.
    :type str: str
    :return: a string that only contains the first two letters.

    """
    if len(str) ==1:
        return str[0]
    return str[0] + str[1]


def generateSpell(sentence, model, oword):
    """
    Generates a Spell from a sentence.

    :param sentence: string which is the definition of the spell you want to create.
    :type sentence: str
    :return: list containing the spell and the spell type.
    :param model: loaded vector orepresentation of words.
        :type model: data file loaded.
    """

    spell = []
    vector,temp_bogus = sentenceToWord(sentence, model, oword)
    vector = vector[0]
    scale = generateScale(count_instances('spell_prob.csv'))
    selection = random.random()
    spell_meta = getSpellType(scale, selection)

    try:
        target_lang = langCode(spell_meta[1])
    except:
        log("langCode function didn't work. Using default latin.")
        target_lang = "la"

    if target_lang == "PL":
        spell.append(pigLatin(vector))
    else:
        spell.append(translate2(vector, target_lang))
    spell.append(spell_meta[0])
    spell.append(vector) #The original word before translation is also added onto the end for evaluation
purposes.
    return spell, temp_bogus


def load_vectors(path, is_binary):
    """
    This loads the vectors supplied by the path.
    @param path: The path to the vector file
    @type path: str
    @param is_binary: states whether file is a binary file.
    @type is_binary: boolean
    """
    print("Loading: ", path)
    model = gensim.models.Word2Vec.load_word2vec_format(path, binary=is_binary)
    model.init_sims(replace=True)
    print("Loaded: ", path)
    return model
```

```python
def is_synonym(n_word, o_word):
    """
    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.
    @param n_word: The new word generated.
    @type n_word: str
    @param o_word: The original word in the definition.
    @type o_word: str
    """
    synonyms=[]
    synsets = wordnet.synsets(o_word)
    for synset in synsets:
        synonyms = synonyms+ synset.lemma_names()

    return n_word in synonyms


def run_experiment(model, num_experiments):
    """
        This function runs the experiments with the paramters set.
        It then returns all the necessary data for processing and output.
        @param model: The vectors loaded.
        @type model: The loaded vector object
        @param num_experiments: The number of experiments to run.
        @type num_experiments: int
    """
    average = 0.0
    iterationCount = 0
    scores = []
    cos_dists = []
    avg_cos_dists = []
    syn_experiments = []
    bword_counts = []
    scores_per_spell=[[] for x in range(10)]#tracks each spell score MUST BE CHANGED TO NUM ENTRIES.
    table1 = []
    table2 = []
    bwords_spell= [[] for x in range(10)] #tracks the number of bogus words against size
    for i in range(0, num_experiments):
        table1 = []
        table2 = []
        print("---------------", i, "---------------")
        log("---------------"+str(i) +  "---------------")
        bogus_words = 0
        spellFile = open("spells.csv")
        entry = []
        score = 0
        count = 0
        syn_counts = 0
        for line in spellFile:
            count+=1
            line = line.strip("\n")
            entry = line.split(",")
            #sen_len.append(len(entry[1].split(" ")))#records length of the sentence.
            #print(len(entry[1].split(" ")))

            spell, temp_bogus = generateSpell(entry[1], model,entry[3] )
            bwords_spell[len(entry[1].split(" "))].append(temp_bogus) #stores the bogus words.
            bogus_words+= temp_bogus
            if args.verbose:
                print("Your new spell is: ", spell[0])
            if spell[2].lower() not in entry[1].split():
                score +=1
                scores_per_spell[len(entry[1].split(" "))].append(1) #keeps track of originality scores.
            else:
                scores_per_spell[len(entry[1].split(" "))].append(0)
            table1.append([spell[0]])
            table2.append([spell[2]])
            #calculate the cosine similarity.
            og_wd = model[entry[-1].strip()]
            nw_wd = model[spell[-1]]
            cos_dists.append(distance.cosine(og_wd, nw_wd))#added log to improve output graph.
            if is_synonym(spell[2].lower(), entry[-1]):
                syn_counts +=1
        #print(tabulate(table1,headers=["Translated"]))
        print("Experiment Results")

        print("Num of spells that don't feature in definition: ", score)
        print("Percentage: ", ((float(score)/count) * 100),"%")
        print("Average Cosine-simalarity:", float(sum(cos_dists) / len(cos_dists)))
        print("Num of spells which are synonyms: ", syn_counts)
        print("Num of words selected that are not real words: ", bogus_words)
        scores.append((float(score)/count) * 100)
        syn_experiments.append(syn_counts)
        bword_counts.append(bogus_words)
        spellFile.close()
        iterationCount +=1
```

```python
            average += (float(score)/count)*100
            avg_cos_dists.append(float(sum(cos_dists) / len(cos_dists)))
    return scores, syn_experiments,average, avg_cos_dists, iterationCount, bword_counts, scores_per_spell,
bwords_spell


# ============================================================================
# Main part of the program.
# ============================================================================
if __name__ == '__main__':
    parser = argparse.ArgumentParser(
            'Use Word2Vec or GloVe datasets to generate Harry Potter Spells')
    parser.add_argument('--glove', action='store_const', const = 'glove',
            help='Use the GloVe dataset instead of the default Word2Vec.')
    parser.add_argument('--exp',
    help="Specifies the number of experiments on this run. Default is 20.",
            action='store', type=int)
    parser.add_argument('--verbose', action='store_const', const = 'verbose',
            help='Prints out the spell names')
    parser.add_argument('--comp', action= 'store_const', const='comp',
            help = "Runs the word2vec vectors, and the GloVe vectors")
    args = parser.parse_args()

    logFile = open("log.txt", 'w' ) #the log file is blank at start of each execution
    logFile.close() #closes the log file
    num_experiments = 20
    if args.exp != None:
        num_experiments = args.exp

    if args.comp: # comparison mode.
        print("Compare Mode")
        log("----------------------Compare Mode----------------------")
        print("Vectors used: Word2Vec")
        log("---------------"+ "Vectors used: Word2Vec"+ "---------------")
        model = load_vectors("../../vectors/GoogleNews-vectors-negative300.bin", True)

        #Run word2vec experiments and then stores data in dataframe.
        w_scores, w_syn_experiments, w_average, w_avg_cos_dists, iterationCount, w_bword_counts,
w_spells_per, w_bwords_per= run_experiment(model, num_experiments)
        w_vec=["word2vec" for x in w_scores]
        del model
        print("Vectors used: GloVe")
        log("---------------" +  "Vectors used: GloVe"+ "---------------")
        model = load_vectors("../../vectors/glove.txt.vw", False)

        # run experiments and move results into data frame.
        g_scores, g_syn_experiments, g_average, g_avg_cos_dists, iterationCount, g_bword_counts,
g_spells_per, g_bwords_per= run_experiment(model, num_experiments)
        g_vec = ["glove" for x in g_scores]

        scores=w_scores + g_scores
        syn_experiments = w_syn_experiments + g_syn_experiments
        avg_cos_dists = w_avg_cos_dists + g_avg_cos_dists
        bword_counts = w_bword_counts + g_bword_counts
        vectors = w_vec + g_vec

        ##for the ts plots
        g_vec = ["GloVe" for x in g_spells_per]
        w_vec = ["Word2Vec" for x in w_spells_per]
        bwords_per = w_bwords_per + g_bwords_per
        spells_per = w_spells_per + g_spells_per
        vec = w_vec + g_vec
        ##adds values for empty rows.#might want to remove empty rows later.
        for row in spells_per:
            if len(row) == 0:
                row.append(0)
        for row in bwords_per:
            if len(row) == 0:
                row.append(0)

        spells_per_avg = [float(sum(l)/len(l)) for l in spells_per]
        length= [x for x in range(1, len(w_spells_per)+1)] + [x for x in range(1, len(g_spells_per)+1)]
        bwords_per_avg = [float(sum(l)/len(l)) for l in bwords_per]
        len_results = pd.DataFrame({"originality":spells_per_avg,"length":length,"bwords":bwords_per_avg,
"vectors":vec})

        box_len = []
        box_score= []
        box_vec=[]


        for i in range(0,len(w_spells_per)):
            for row in w_spells_per[i]:
                box_len.append(i+1)
                box_score.append(row)
                box_vec.append("word2vec")
```

```python
415
416                 for row2 in g_spells_per[i]:
417                     box_len.append(i+1)
418                     box_score.append(row2)
419                     box_vec.append("GloVe")
420
421         box_data = pd.DataFrame({"length":box_len, "originality":box_score, "vectors":box_vec})
            #originality vs size plots.
422         ax = sns.tsplot(time="length", value="originality",
    unit="vectors",condition="vectors",data=len_results   )
423 #        sns.plt.xticks([0,1,2,3,4,5,6,7,8,9,10])
424         sns.plt.show()
425
426         ax = sns.distplot(box_score)
427         sns.plt.show()
            #box plot
428         ax = sns.boxplot(x="length", y = "originality", hue="vectors", data=box_data)
429         sns.plt.show()
430
431         box_len = []
432         box_score= []
            box_vec=[]
433
434         for i in range(0,len(w_bwords_per)):
435             for row in w_bwords_per[i]:
436                 box_len.append(i+1)
437                 box_score.append(row)
                    box_vec.append("word2vec")
438             for row2 in g_bwords_per[i]:
439                 box_len.append(i+1)
440                 box_score.append(row2)
441                 box_vec.append("GloVe")
442         box_data = pd.DataFrame({"length":box_len, "bwords":box_score, "vectors":box_vec})
443         #gibberish vs size plots
444         ax = sns.tsplot(time="length", value="bwords", unit="vectors",condition="vectors",data=len_results
    )
445 #        sns.plt.xticks([0,1,2,3,4,5,6,7,8,9,10])
446         sns.plt.show()
447         #histogram
448         ax = sns.distplot(box_score)
449         sns.plt.show()
            #box plot
450         ax = sns.boxplot(x="length", y = "bwords", hue="vectors", data=box_data)
451         #sns.despine(offset=10, trim=True)
452         sns.plt.show()
            #
453
454
455         ##output results.
456
457         print("----------------word2vec Experiment Results------------------")
458         print("The mean average percentage over ", iterationCount , "tests: ",
                 (w_average/iterationCount), "%")
459         print("The mean cosine simalarity over ", iterationCount, "tests: ",
460                 float(sum(w_avg_cos_dists)/ len(w_avg_cos_dists)))
            print("The mean amount of synonyms", (sum(w_syn_experiments)/ iterationCount))
461         print("Average number of words that are not fit for translation:
462 ",float(sum(w_bword_counts)/iterationCount))
463
464         print("----------------GloVe Experiment Results------------------")
465         print("The mean average percentage over ", iterationCount , "tests: ",
466                 (g_average/iterationCount), "%")
467         print("The mean cosine simalarity over ", iterationCount, "tests: ",
468                 float(sum(g_avg_cos_dists)/ len(g_avg_cos_dists)))
            print("The mean amount of synonyms", (sum(g_syn_experiments)/ iterationCount))
469         print("Average number of words that are not fit for translation:
470 ",float(sum(g_bword_counts)/iterationCount))
471
472
473         results = pd.DataFrame({"scores":scores, "similarity":avg_cos_dists, "synonyms":syn_experiments,
    "vectors":vectors, "bwords":bword_counts})
474
475
476         sim = sns.violinplot(x="vectors", y="similarity", data=results)
477         sns.plt.title("Comparison of Similarity over "+str( iterationCount)+ " experiments")
478         sns.plt.show()
            sc = sns.violinplot(x="vectors", y="scores", data=results)
479         sns.plt.title("Comparison of accuracy scores over "+str(iterationCount)+ " experiments")
480         sns.plt.show()
            bw = sns.violinplot(x="vectors", y="bwords", data=results)
481         sns.plt.title("Comparison of invalid words over "+ str(iterationCount)+ " experiments")
482         sns.plt.show()
483         sns.plt.title("Comparison of synonyms over " +str( iterationCount) +" experiments")
```

```python
        syn = sns.violinplot(x="vectors", y="synonyms", data=results)
        sns.plt.show()

    else: # test an individual mode.
        if args.glove:
            print("Vectors used: GloVe")
            log("--------------" +  "Vectors used: GloVe"+ "--------------")
            model = load_vectors("../../vectors/glove.txt.vw", False)
        else:
            print("Vectors used: Word2Vec")
            log("--------------"+ "Vectors used: Word2Vec"+ "--------------")
            model = load_vectors("../../vectors/GoogleNews-vectors-negative300.bin", True)


        scores, syn_experiments, average, avg_cos_dists, iterationCount, bword_counts,  spells_per=
run_experiment(model, num_experiments)
        print("---------------Experiment Results------------------")
        print("The mean average percentage over ", iterationCount , "tests: ",
                (average/iterationCount), "%")
        print("The mean cosine simalarity over ", iterationCount, "tests: ",
                float(sum(avg_cos_dists)/ len(avg_cos_dists)))
        print("The mean amount of synonyms", (sum(syn_experiments)/ iterationCount))
        print("Average number of words that are not fit for translation:
",float(sum(bword_counts)/iterationCount))
        results = pd.DataFrame({'scores': scores, 'similarity': avg_cos_dists})
        #loop through and add an entry to any empty fields.
        for row in spells_per:
            if len(row) == 0:
                row.append(0)

        spells_per_avg = [float(sum(l)/len(l)) for l in spells_per]
        length= [x for x in range(0, len(spells_per_avg))]

        vec = ["vector" for x in spells_per_avg]

        len_results = pd.DataFrame({"scores":spells_per_avg,"length":length, "vec":vec})

        ax = sns.tsplot(time="length", value="scores", unit="vec",condition="vec",data=len_results  )
        sns.plt.show()
       # ts_plot(len_results, "scores")
        ax2 = sns.violinplot(x=results["similarity"])
        sns.plt.show()
        ax = sns.violinplot(x="scores", y="similarity", data=results)
        sns.plt.show()
```