```
Assignment 2 - Introduction to NLTK
          In part 1 of this assignment you will use nltk to explore the Herman Melville novel Moby Dick. Then in part 2 you will create a spelling recommender function
          that uses nltk to find words similar to the misspelling.
          Part 1 - Analyzing Moby Dick
In [13]: #nltk.download()
          import nltk
          import pandas as pd
          import numpy as np
          nltk.download('gutenberg')
          nltk.download('genesis')
          nltk.download('inaugural')
          nltk.download('nps_chat')
          nltk.download('treebank')
          nltk.download('udhr')
          nltk.download('webtext')
          nltk.download('wordnet')
          nltk.download('punkt')
          nltk.download('wordnet')
          #nltk.download()
          from nltk.book import *
          # If you would like to work with the raw text you can use 'moby_raw'
          with open('moby.txt', 'r') as f:
              moby_raw = f.read()
          # If you would like to work with the novel in nltk. Text format you can use 'text1'
          moby_tokens = nltk.word_tokenize(moby_raw)
          text1 = nltk.Text(moby_tokens)
          [nltk data] Downloading package gutenberg to /home/jovyan/nltk data...
          [nltk_data] Package gutenberg is already up-to-date!
          [nltk_data] Downloading package genesis to /home/jovyan/nltk_data...
          [nltk_data] Package genesis is already up-to-date!
          [nltk data] Downloading package inaugural to /home/jovyan/nltk data...
          [nltk_data] Package inaugural is already up-to-date!
          [nltk_data] Downloading package nps_chat to /home/jovyan/nltk_data...
          [nltk_data] Package nps_chat is already up-to-date!
          [nltk_data] Downloading package treebank to /home/jovyan/nltk_data...
          [nltk_data] Package treebank is already up-to-date!
          [nltk_data] Downloading package udhr to /home/jovyan/nltk_data...
          [nltk_data] Unzipping corpora/udhr.zip.
          [nltk_data] Downloading package webtext to /home/jovyan/nltk_data...
          [nltk_data] Package webtext is already up-to-date!
          [nltk_data] Downloading package wordnet to /home/jovyan/nltk_data...
          [nltk_data] Unzipping corpora/wordnet.zip.
          [nltk_data] Downloading package punkt to /home/jovyan/nltk_data...
          [nltk_data] Unzipping tokenizers/punkt.zip.
          [nltk_data] Downloading package wordnet to /home/jovyan/nltk_data...
          [nltk_data] Package wordnet is already up-to-date!
          Example 1
          How many tokens (words and punctuation symbols) are in text1?
          This function should return an integer.
In [14]: def example_one():
              return len(nltk.word_tokenize(moby_raw)) # or alternatively len(text1)
          example_one()
Out[14]: 254989
          Example 2
          How many unique tokens (unique words and punctuation) does text1 have?
          This function should return an integer.
In [15]: def example_two():
              return len(set(nltk.word tokenize(moby raw))) # or alternatively len(set(text1))
          example_two()
Out[15]: 20755
          Example 3
          After lemmatizing the verbs, how many unique tokens does text1 have?
          This function should return an integer.
In [16]: from nltk.stem import WordNetLemmatizer
          def example three():
              lemmatizer = WordNetLemmatizer()
              lemmatized = [lemmatizer.lemmatize(w,'v') for w in text1]
              return len(set(lemmatized))
          example_three()
Out[16]: 16900
          Question 1
          What is the lexical diversity of the given text input? (i.e. ratio of unique tokens to the total number of tokens)
          This function should return a float.
In [17]: def answer_one():
              return example_two()/example_one()
          answer_one()
Out[17]: 0.08139566804842562
          Question 2
          What percentage of tokens is 'whale'or 'Whale'?
          This function should return a float.
In [18]: def answer_two():
              return (text1.vocab()['whale'] + text1.vocab()['Whale']) / len(nltk.word_tokenize(moby_raw)) * 100
          answer_two()
Out[18]: 0.4125668166077752
          Question 3
          What are the 20 most frequently occurring (unique) tokens in the text? What is their frequency?
          This function should return a list of 20 tuples where each tuple is of the form (token, frequency). The list should be sorted in descending order of
          frequency.
In [20]: def answer_three():
              import operator
              return sorted(text1.vocab().items(), key=operator.itemgetter(1), reverse=True)[:20]
          answer_three()
Out[20]: [(',', 19204),
           ('the', 13715),
           ('.', 7308),
           ('of', 6513),
           ('and', 6010),
           ('a', 4545),
           ('to', 4515),
           (';', 4173),
           ('in', 3908),
           ('that', 2978),
           ('his', 2459),
           ('it', 2196),
           ('I', 2097),
           ('!', 1767),
           ('is', 1722),
           ('--', 1713),
           ('with', 1659),
           ('he', 1658),
           ('was', 1639),
           ('as', 1620)]
          Question 4
          What tokens have a length of greater than 5 and frequency of more than 150?
          This function should return an alphabetically sorted list of the tokens that match the above constraints. To sort your list, use sorted()
In [21]: def answer_four():
              return sorted([token for token, freq in text1.vocab().items() if len(token) > 5 and freq > 150])
          answer_four()
Out[21]: ['Captain',
           'Pequod',
           'Queequeg',
           'Starbuck',
           'almost',
           'before',
           'himself',
           'little',
           'seemed',
           'should',
           'though',
           'through',
           'whales',
           'without']
          Question 5
          Find the longest word in text1 and that word's length.
          This function should return a tuple (longest word, length).
In [22]: def answer_five():
              import operator
              return sorted([(token, len(token))for token, freq in text1.vocab().items()], key=operator.itemgetter(1), reverse=T
          rue)[0]
          answer_five()
Out[22]: ("twelve-o'clock-at-night", 23)
          Question 6
          What unique words have a frequency of more than 2000? What is their frequency?
          "Hint: you may want to use isalpha() to check if the token is a word and not punctuation."
          This function should return a list of tuples of the form (frequency, word) sorted in descending order of frequency.
In [23]: def answer_six():
              import operator
              return sorted([(freq, token) for token, freq in text1.vocab().items() if freq > 2000 and token.isalpha()], key=ope
          rator.itemgetter(0), reverse=True)
          answer_six()
Out[23]: [(13715, 'the'),
           (6513, 'of'),
           (6010, 'and'),
           (4545, 'a'),
           (4515, 'to'),
           (3908, 'in'),
           (2978, 'that'),
           (2459, 'his'),
           (2196, 'it'),
           (2097, 'I')]
          Question 7
          What is the average number of tokens per sentence?
          This function should return a float.
In [30]: def answer_seven():
              # use the built-in package to split the text into sentences
              sentences = nltk.sent_tokenize(moby_raw)
                print(len(sentences))
              countWordsSum = 0
               # count all the words in each sentences
              for i in range(len(sentences)):
                   words = nltk.word_tokenize(sentences[i])
                   countWordsSum = countWordsSum + len(words)
              return (countWordsSum / len(sentences))
          answer_seven()
Out[30]: 25.881952902963864
          Question 8
          What are the 5 most frequent parts of speech in this text? What is their frequency?
          This function should return a list of tuples of the form (part_of_speech, frequency) sorted in descending order of frequency.
In [38]: moby_frequencies = nltk.FreqDist(moby_tokens)
               # set up the dataframe
          df = pd.DataFrame(moby_frequencies.most_common(),columns=["token", "frequency"])
          # find the valid words in moby
          moby_words = df[df.token.str.isalpha()]
          print(moby_words)
                           token frequency
                                      13715
          1
                             the
          3
                              of
                                       6513
                                        6010
                             and
          5
                               a
                                       4545
                                        4515
                              to
                              in
                                       3908
          9
                                       2978
                            that
          10
                             his
                                       2459
          11
                              it
                                       2196
          12
                              I
                                       2097
          14
                              is
                                       1722
                                       1659
          16
                           with
          17
                                       1658
                              he
          18
                                       1639
                             was
          19
                                       1620
                              as
          23
                             all
                                       1444
          24
                             for
                                       1413
          25
                                       1280
                            this
          26
                              at
                                       1230
          27
                                       1170
                             not
          28
                                       1135
                              by
          29
                             but
                                       1110
          30
                             him
                                       1058
          31
                            from
                                       1052
          32
                                       1027
                              be
          34
                                       1003
                              on
          35
                                         914
                              so
          36
                                         880
                             one
          37
                                         841
                            you
          38
                                         782
                          whale
          • • •
                                         . . .
          20720
                       lookouts
                                          1
          20722
                        animate
          20723
                      inanimate
          20724
                      whelmings
          20725
                  intermixingly
          20726
                       ironical
          20727
                    coincidings
          20728
                     destroying
                     backwardly
          20729
          20731
                     tauntingly
          20732
                    incommoding
          20733
                      intercept
          20734
                       etherial
          20735
                         thrill
                                           1
          20737
                       Epilogue
          20738
                           ONLY
          20739
                        ESCAPED
          20740
                           THEE
          20741
                      halfspent
          20742
                        suction
          20743
                        closing
          20745
                          Ixion
                           Till
          20746
          20747
                      liberated
          20748
                         Buoyed
                      dirgelike
          20749
                       padlocks
          20750
          20751
                       sheathed
          20753
                      retracing
          20754
                         orphan
          [18464 rows x 2 columns]
In [42]: def answer_eight():
               import collections
              import nltk
              nltk.download('averaged_perceptron_tagger')
              from collections import Counter
              import operator
              return sorted(Counter([tag for token, tag in nltk.pos_tag(text1)]).items(), key=operator.itemgetter(1), reverse=Tr
          ue)[:5]
          answer_eight()
          [nltk_data] Downloading package averaged_perceptron_tagger to
                           /home/jovyan/nltk_data...
          [nltk data]
          [nltk_data]
                         Package averaged_perceptron_tagger is already up-to-
          [nltk_data]
Out[42]: [('NN', 32730), ('IN', 28657), ('DT', 25867), (',', 19204), ('JJ', 17620)]
          Part 2 - Spelling Recommender
          For this part of the assignment you will create three different spelling recommenders, that each take a list of misspelled words and recommends a correctly
          spelled word for every word in the list.
          For every misspelled word, the recommender should find the word in correct spellings that has the shortest distance*, and starts with the same letter
          as the misspelled word, and return that word as a recommendation.
          *Each of the three different recommenders will use a different distance measure (outlined below).
          Each of the recommenders should provide recommendations for the three default words provided: ['cormulent', 'incendence', 'validrate'].
In [44]: from nltk.corpus import words
          import nltk
          nltk.download('words')
          correct_spellings = words.words()
          [nltk_data] Downloading package words to /home/jovyan/nltk_data...
          [nltk data] Unzipping corpora/words.zip.
          Question 9
          For this recommender, your function should provide recommendations for the three default words provided above using the following distance metric:
          <u>Jaccard distance</u> on the trigrams of the two words.
          This function should return a list of length three: ['cormulent_reccomendation', 'incendence_reccomendation',
           'validrate_reccomendation'].
In [45]: def answer_nine(entries=['cormulent', 'incendenece', 'validrate']):
              result = []
              import operator
              for entry in entries:
                   spell_list = [spell for spell in correct_spellings if spell.startswith(entry[0]) and len(spell) > 2]
                   distance_list = [(spell, nltk.jaccard_distance(set(nltk.ngrams(entry, n=3)), set(nltk.ngrams(spell, n=3)))) fo
          r spell in spell list]
                   result.append(sorted(distance_list, key=operator.itemgetter(1))[0][0])
              return result # Your answer here
          answer_nine()
Out[45]: ['corpulent', 'indecence', 'validate']
          Question 10
          For this recommender, your function should provide recommendations for the three default words provided above using the following distance metric:
          Jaccard distance on the 4-grams of the two words.
          This function should return a list of length three: ['cormulent_reccomendation', 'incendence_reccomendation',
           'validrate reccomendation'].
In [46]: def answer_ten(entries=['cormulent', 'incendenece', 'validrate']):
              result = []
              import operator
              for entry in entries:
                   spell_list = [spell for spell in correct_spellings if spell.startswith(entry[0]) and len(spell) > 2]
                   distance_list = [(spell, nltk.jaccard_distance(set(nltk.ngrams(entry, n=4)), set(nltk.ngrams(spell, n=4)))) fo
          r spell in spell_list]
                   result.append(sorted(distance_list, key=operator.itemgetter(1))[0][0])
              return result # Your answer here
          answer_ten()
Out[46]: ['cormus', 'incendiary', 'valid']
          Question 11
          For this recommender, your function should provide recommendations for the three default words provided above using the following distance metric:
          Edit distance on the two words with transpositions.
          This function should return a list of length three: ['cormulent reccomendation', 'incendence reccomendation',
           'validrate_reccomendation'].
```

In [47]: def answer_eleven(entries=['cormulent', 'incendenece', 'validrate']):

result.append(sorted(distance_list, key=operator.itemgetter(1))[0][0])

spell_list = [spell for spell in correct_spellings if spell.startswith(entry[0]) and len(spell) > 2]

distance_list = [(spell, nltk.edit_distance(entry, spell, transpositions=True)) for spell in spell_list]

result = []

answer_eleven()

In []:

import operator

for entry in entries:

Out[47]: ['corpulent', 'intendence', 'validate']

return result# Your answer here

You are currently looking at **version 1.0** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera

platform, visit the Jupyter Notebook FAQ course resource.