

```
In [1]: import numpy as np
import sys
import importlib
import nltk
nltk.download('punkt')

(nltk_data) Downloading package punkt to
(nltk_data) /Users/adilhoulam/nltk_data...
(nltk_data) Package punkt is already up-to-date!

Out[1]: True

In [3]: #Inspired from the context2vec evaluation function explore_context2vec
import sys
import numpy
import six
import sys
import traceback
import re

from chainer import cuda
from context2vec.common.context_models import Toks
from context2vec.common.model_reader import ModelReader

class ParseException(Exception):
    def __init__(self, str):
        super(ParseException, self).__init__(str)

target_exp = re.compile('\[.*\]')

def parse_index(line):
    sent = line.strip().split()
    target_pos = None
    for i, word in enumerate(sent):
        if target_exp.match(word) != None:
            target_pos = i
            if word == '[':
                word = None
            else:
                word = word[i+1:]
            sent[i] = word
    return sent, target_pos

def multi_sim(w, target_v, context_v):
    target_similarity = w.dot(target_v)
    target_similarity[target_similarity<0] = 0.0
    context_similarity = w.dot(context_v)
    context_similarity[context_similarity<0] = 0.0
    return (target_similarity * context_similarity)

n_rows = 1000 # number of search result to show
gpu = -1 # todo: make this work with gpu

if gpu >= 0:
    cuda.check_cuda_available()
    cuda.get_device_name(gpu, use())
    np = cuda.cupy if gpu >= 0 else numpy

model_reader = ModelReader('/Users/adilhoulam/Downloads/context2vec.ukwac.model.package.\
backslash/context2vec.ukwac.model.params')
w = model_reader.w
vocab_index = model_reader.word_index
index2word = model_reader.index2word
model = model_reader.model

def find_most_similar(line, n_result):
    sent, target_pos = parse_input(line)
    if target_pos == None:
        raise ParseException("Can't find the target position.")
    if sent[target_pos] == None:
        target_v = None
    elif sent[target_pos] not in word2index:
        raise ParseException("Target word is out of vocabulary.")
    else:
        target_v = w[word2index[sent[target_pos]]]
        if len(sent) > 1:
            context_v = model.context2vec(sent, target_pos)
            context_v = context_v / np.sqrt(context_v * context_v).sum())
        else:
            context_v = None

        if target_v is not None and context_v is not None:
            similarity = multi_sim(w, target_v, context_v)
        else:
            if target_v is not None:
                v = target_v
            elif context_v is not None:
                v = context_v
            else:
                raise ParseException("Can't find a target nor context.")
            similarity = (v.dot(v)+1.0)/2 # Cosine similarity can be negative, mapping similar
            ty co

    resultat = []
    for i in (-similarity).argsort():
        if i < 0:
            continue
        resultat.append([index2word[i],similarity[i]])
        print('({})'.format(index2word[i],similarity[i]))
        count += 1
        if count == n_result:
            break
    return resultat

/anaconda2/lib/python2.7/site-packages/h5py/_init_.py:36: FutureWarning: Conversion of the
second argument of issubdtype from float to np.floating is deprecated. In future, it will
be converted to np.float64 == np.dtype(float).type
from _conv import register_converters as _register_converters

Reading config file: /Users/adilhoulam/Downloads/context2vec.ukwac.model.package/context2ve
c.ukwac.model.params
Config: {'config_path': '/Users/adilhoulam/Downloads/context2vec.ukwac.model.package/', 'm
odel_file': 'context2vec.ukwac.model', 'deep': 'yes', 'drop_ratio': '0.0', 'words_file': 'co
ntext2vec.ukwac.words.targets', 'unit': '300'}

In [10]: import re
from collections import Counter

words = re.findall(r'w+', open('/Users/adilhoulam/Downloads/words.txt').read()).lower()

Dict = Counter(words)

def invert(words):
    return set(w for w in words if w in Dict)

def candidates(word):
    return (known[word]) or known(edit1(word)) or known(edit2(word)) or (word)

def f(words, N=sum(Dict.values())):
    return Dict[word] / N

def edit1(word):
    "All edits that are one edit away from 'word'."
    letters = 'abcdefghijklmnopqrstuvwxyz'
    splits = [(word[i], word[i+1]) for i in range(len(word)-1)]
    delchars = [(i, j) for i, j in splits if word[i] != word[j]]
    transposes = [(i, j) for i, j in splits if word[i] != word[j]]
    replaces = [(i, c) for i, c in letters if word[i] != c]
    inserts = [(i, c) for i, c in letters if word[i] != c]
    return set(delchars + transposes + replaces + inserts)

def edit2(word):
    "All edits that are two edits away from 'word'."
    return set(e2 for e1 in edit1(word) for e2 in edit1(e1))

def edit_distance(s1,s2):
    n = np.zeros((len(s1),len(s2)),dtype=int)
    m[0,0] = np.arange(len(s1))
    m[0,1] = np.arange(len(s2))
    for i in range(1,len(s1)):
        for j in range(1,len(s2)):
            if s1[i-1]==s2[j-1]:
                m[i,j] = min(m[i-1,j],m[i,j-1],m[i-1,j-1])
            else:
                m[i,j] = min(m[i-1,j]+1,m[i,j-1]+1,m[i-1,j-1]+1)
    return m[len(s1)-1,len(s2)-1]
    return m[len(s1)-1,len(s2)-1]

def correction_context(word,candidates):
    edit1 = []
    edit2 = []
    for i in range(len(candidates)):
        dist = edit_distance(word,candidates[i][0])
        if dist<1:
            edit1.append(candidates[i][0])
        if dist==2:
            edit2.append(candidates[i][0])
    if len(edit1)==0:
        return edit2[0]
    return correction_spelling(word)
    return edit1[0]

def correction_spelling(word):
    return max(candidates(word), key=P)

In [11]: path = '/Users/adilhoulam/Downloads/CorpusBataclan_en.IM.raw.txt'
def load(path,samples):
    with open(path) as f:
        lines = f.readlines()
        sentences = []
        for i in range(samples):
            sentences.append(" ".join(filter(lambda x: x[0]!='&' and x[0]!='&' and \
sentences.append(" ".join(filter(lambda x: x[0]!='&' and x[0]!='&' and \
sentences.append(" ".join(filter(lambda x: x[0]!='&' and x[0]!='&' and \
tokens_1 = []
for i in range(len(sentences)):
    tokens_1.append(nltk.word_tokenize(sentences[i]))

tokens_1,nouns_1,nouns_2 = parsefile(path,30000)

In [12]: def text_to_correct(nouns_1,nouns_2,tokens_1):
    corrected = nouns_2
    for i in range(len(nouns_1)):
        if (100==0):
            print('='*10)
            for p in range(len(nouns_1[i])):
                if nouns_1[i][p] in Dict and len(nouns_1[i][p])>1:
                    corrected[i][p] = nouns_1[i][p]
            continue
            elif len(nouns_1[i][p])>0:
                corrected[i][p] = tokens_1[i][p]
            else:
                word = nouns_1[i][p]
                line = ''
                for p in range(0,j):
                    if p==i:
                        line = line + nouns_1[i][p]
                    else:
                        line = line + ' ' + nouns_1[i][p]
                if len(line)==0:
                    line = []
                else:
                    line = line + ' '
                for p in range(1,len(nouns_1[i])):
                    line=line+ ' '+nouns_1[i][p]
                    if " ".join(line.split())!=[]:
                        corrected[i][p] = correction_spelling(word)
                else:
                    candidates = find_most_similar(line,10)
                    corrected[i][p] = correction_context(word,candidates_)
                return corrected

corrected = text_to_correct(nouns_1,nouns_2)

('1', 0)
('1', 1)
('1', 200)
('1', 300)
('1', 400)
('1', 500)
('1', 600)
('1', 700)
('1', 800)
('1', 900)
('1', 1000)
('1', 1100)
('1', 1200)
('1', 1300)
('1', 1400)
('1', 1500)
('1', 1600)
('1', 1700)
('1', 1800)
('1', 1900)
('1', 2000)
('1', 2100)
('1', 2200)
('1', 2300)
('1', 2400)
('1', 2500)
('1', 2600)
('1', 2700)
('1', 2800)
('1', 2900)
('1', 3000)
('1', 3100)
('1', 3200)
('1', 3300)
('1', 3400)
('1', 3500)
('1', 3600)
('1', 3700)
('1', 3800)
('1', 3900)
('1', 4000)
('1', 4100)
('1', 4200)
('1', 4300)
('1', 4400)
('1', 4500)
('1', 4600)
('1', 4700)
('1', 4800)
('1', 4900)
('1', 5000)
('1', 5100)
('1', 5200)
('1', 5300)
('1', 5400)
('1', 5500)
('1', 5600)
('1', 5700)
('1', 5800)
('1', 5900)
('1', 6000)
('1', 6100)
('1', 6200)
('1', 6300)
('1', 6400)
('1', 6500)
('1', 6600)
('1', 6700)
('1', 6800)
('1', 6900)
('1', 7000)
('1', 7100)
('1', 7200)
('1', 7300)
('1', 7400)
('1', 7500)
('1', 7600)
('1', 7700)
('1', 7800)
('1', 7900)
('1', 8000)
('1', 8100)
('1', 8200)
('1', 8300)
('1', 8400)
('1', 8500)
('1', 8600)
('1', 8700)
('1', 8800)
('1', 8900)
('1', 9000)
('1', 9100)
('1', 9200)
('1', 9300)
('1', 9400)
('1', 9500)
('1', 9600)
('1', 9700)
('1', 9800)
('1', 9900)
('1', 10000)
('1', 10100)
('1', 10200)
('1', 10300)
('1', 10400)
('1', 10500)
('1', 10600)
('1', 10700)
('1', 10800)
('1', 10900)
('1', 11000)
('1', 11100)
('1', 11200)
('1', 11300)
('1', 11400)
('1', 11500)
('1', 11600)
('1', 11700)
('1', 11800)
('1', 11900)
('1', 12000)
('1', 12100)
('1', 12200)
('1', 12300)
('1', 12400)
('1', 12500)
('1', 12600)
('1', 12700)
('1', 12800)
('1', 12900)
('1', 13000)
('1', 13100)
('1', 13200)
('1', 13300)
('1', 13400)
('1', 13500)
('1', 13600)
('1', 13700)
('1', 13800)
('1', 13900)
('1', 14000)
('1', 14100)
('1', 14200)
('1', 14300)
('1', 14400)
('1', 14500)
('1', 14600)
('1', 14700)
('1', 14800)
('1', 14900)
('1', 15000)
('1', 15100)
('1', 15200)
('1', 15300)
('1', 15400)
('1', 15500)
('1', 15600)
('1', 15700)
('1', 15800)
('1', 15900)
('1', 16000)
('1', 16100)
('1', 16200)
('1', 16300)
('1', 16400)
('1', 16500)
('1', 16600)
('1', 16700)
('1', 16800)
('1', 16900)
('1', 17000)
('1', 17100)
('1', 17200)
('1', 17300)
('1', 17400)
('1', 17500)
('1', 17600)
('1', 17700)
('1', 17800)
('1', 17900)
('1', 18000)
('1', 18100)
('1', 18200)
('1', 18300)
('1', 18400)
('1', 18500)
('1', 18600)
('1', 18700)
('1', 18800)
('1', 18900)
('1', 19000)
('1', 19100)
('1', 19200)
('1', 19300)
('1', 19400)
('1', 19500)
('1', 19600)
('1', 19700)
('1', 19800)
('1', 19900)
('1', 20000)
('1', 20100)
('1', 20200)
('1', 20300)
('1', 20400)
('1', 20500)
('1', 20600)
('1', 20700)
('1', 20800)
('1', 20900)
('1', 21000)
('1', 21100)
('1', 21200)
('1', 21300)
('1', 21400)
('1', 21500)
('1', 21600)
('1', 21700)
('1', 21800)
('1', 21900)
('1', 22000)
('1', 22100)
('1', 22200)
('1', 22300)
('1', 22400)
('1', 22500)
('1', 22600)
('1', 22700)
('1', 22800)
('1', 22900)
('1', 23000)
('1', 23100)
('1', 23200)
('1', 23300)
('1', 23400)
('1', 23500)
('1', 23600)
('1', 23700)
('1', 23800)
('1', 23900)
('1', 24000)
('1', 24100)
('1', 24200)
('1', 24300)
('1', 24400)
('1', 24500)
('1', 24600)
('1', 24700)
('1', 24800)
('1', 24900)
('1', 25000)
('1', 25100)
('1', 25200)
('1', 25300)
('1', 25400)
('1', 25500)
('1', 25600)
('1', 25700)
('1', 25800)
('1', 25900)
('1', 26000)
('1', 26100)
('1', 26200)
('1', 26300)
('1', 26400)
('1', 26500)
('1', 26600)
('1', 26700)
('1', 26800)
('1', 26900)
('1', 27000)
('1', 27100)
('1', 27200)
('1', 27300)
('1', 27400)
('1', 27500)
('1', 27600)
('1', 27700)
('1', 27800)
('1', 27900)
('1', 28000)
('1', 28100)
('1', 28200)
('1', 28300)
('1', 28400)
('1', 28500)
('1', 28600)
('1', 28700)
('1', 28800)
('1', 28900)
('1', 29000)
('1', 29100)
('1', 29200)
('1', 29300)
('1', 29400)
('1', 29500)
('1', 29600)
('1', 29700)
('1', 29800)
('1', 29900)
('1', 30000)
('1', 30100)
('1', 30200)
('1', 30300)
('1', 30400)
('1', 30500)
('1', 30600)
('1', 30700)
('1', 30800)
('1', 30900)
('1', 31000)
('1', 31100)
('1', 31200)
('1', 31300)
('1', 31400)
('1', 31500)
('1', 31600)
('1', 31700)
('1', 31800)
('1', 31900)
('1', 32000)
('1', 32100)
('1', 32200)
('1', 32300)
('1', 32400)
('1', 32500)
('1', 32600)
('1', 32700)
('1', 32800)
('1', 32900)
('1', 33000)
('1', 33100)
('1', 33200)
('1', 33300)
('1', 33400)
('1', 33500)
('1', 33600)
('1', 33700)
('1', 33800)
('1', 33900)
('1', 34000)
('1', 34100)
('1', 34200)
('1', 34300)
('1', 34400)
('1', 34500)
('1', 34600)
('1', 34700)
('1', 34800)
('1', 34900)
('1', 35000)
('1', 35100)
('1', 35200)
('1', 35300)
('1', 35400)
('1', 35500)
('1', 35600)
('1', 35700)
('1', 35800)
('1', 35900)
('1', 36000)
('1', 36100)
('1', 36200)
('1', 36300)
('1', 36400)
('1', 36500)
('1', 36600)
('1', 36700)
('1', 36800)
('1', 36900)
('1', 37000)
('1', 37100)
('1', 37200)
('1', 37300)
('1', 37400)
('1', 37500)
('1', 37600)
('1', 37700)
('1', 37800)
('1', 37900)
('1', 38000)
('1', 38100)
('1', 38200)
('1', 38300)
('1', 38400)
('1', 38500)
('1', 38600)
('1', 38700)
('1', 38800)
('1', 38900)
('1', 39000)
('1', 39100)
('1', 39200)
('1', 39300)
('1', 39400)
('1', 39500)
('1', 39600)
('1', 39700)
('1', 39800)
('1', 39900)
('1', 40000)
('1', 40100)
('1', 40200)
('1', 40300)
('1', 40400)
('1', 40500)
('1', 40600)
('1', 40700)
('1', 40800)
('1', 40900)
('1', 41000)
('1', 41100)
('1', 41200)
('1', 41300)
('1', 41400)
('1', 41500)
('1', 41600)
('1', 41700)
('1', 41800)
('1', 41900)
('1', 42000)
('1', 42100)
('1', 42200)
('1', 42300)
('1', 42400)
('1', 42500)
('1', 42600)
('1', 42700)
('1', 42800)
('1', 42900)
('1', 43000)
('1', 43100)
('1', 43200)
('1', 43300)
('1', 43400)
('1', 43500)
('1', 43600)
('1', 43700)
('1', 43800)
('1', 43900)
('1', 44000)
('1', 44100)
('1', 44200)
('1', 44300)
('1', 44400)
('1', 44500)
('1', 44600)
('1', 44700)
('1', 44800)
('1', 44900)
('1', 45000)
('1', 45100)
('1', 45200)
('1', 45300)
('1', 45400)
('1', 45500)
('1', 45600)
('1', 45700)
('1', 45800)
('1', 45900)
('1', 46000)
('1', 46100)
('1', 46200)
('1', 46300)
('1', 46400)
('1', 46500)
('1', 46600)
('1', 46700)
('1', 46800)
('1', 46900)
('1', 47000)
('1', 47100)
('1', 47200)
('1', 47300)
('1', 47400)
('1', 47500)
('1', 47600)
('1', 47700)
('1', 47800)
('1', 47900)
('1', 48000)
('1', 48100)
('1', 48200)
('1', 48300)
('1', 48400)
('1', 48500)
('1', 48600)
('1', 48700)
('1', 48800)
('1', 48900)
('1', 49000)
('1', 49100)
('1', 49200)
('1', 49300)
('1', 49400)
('1', 49500)
('1', 49600)
('1', 49700)
('1', 49800)
('1', 49900)

In [18]: path = '/Users/adilhoulam/Desktop/test.txt'
def join_text(tokens,path):
    thefile = open(path, 'w')
    lines = []
    for i in range(len(tokens)):
        lines[i].append(" ".join(tokens[i]))
        thefile.write("%s\n" % lines[i])
    return lines

lines = join_text(corrected,path)
```