```python
import numpy as np
import matplotlib.pyplot as plt
from sklearn.metrics import pairwise
from sklearn.manifold import TSNE


defining_model = dict()
f = open('/content/sample_data/vectors.txt', encoding="utf8")

for line in f:
    values = line.split()
    word = values[0]
    coefficient = np.asarray(values[1:], dtype='float32')
    defining_model[word] = coefficient
f.close()
print('Loaded %s word vectors.' % (len(defining_model)))
```

    Loaded 400001 word vectors.

```python
def similar_words(defining_model, word,target_list , num) :
    cosine_dict ={}
    word_list = []
    if(word in defining_model.keys()):
      a = defining_model[word]
    else:
      print('No embeddings found for %s' % word)
      return

    for i in target_list :
        if i != word :
            a = defining_model [i]
            cos_sim = pairwise.cosine_similarity(a.reshape(1, -1),a.reshape(1, -1))
            cosine_dict[i] = cos_sim
    dist_sort=sorted(cosine_dict.items(), key=lambda dist: dist[1],reverse = True)
    for i in dist_sort:
        word_list.append((i[0], i[1]))
    return word_list[0:num]


similar_words(defining_model,'life',list(defining_model.keys()),20)
```

```
[('difference', array([[1.0000004]], dtype=float32)),
 ('forgotten', array([[1.0000004]], dtype=float32)),
 ('chester', array([[1.0000004]], dtype=float32)),
 ('protagonist', array([[1.0000004]], dtype=float32)),
 ('amounting', array([[1.0000004]], dtype=float32)),
 ('laude', array([[1.0000004]], dtype=float32)),
 ('congressionally', array([[1.0000004]], dtype=float32)),
 ('twinkling', array([[1.0000004]], dtype=float32)),
 ('dianna', array([[1.0000004]], dtype=float32)),
 ('replaceable', array([[1.0000004]], dtype=float32)),
 ('cnn.com', array([[1.0000004]], dtype=float32)),
 ('acupuncturist', array([[1.0000004]], dtype=float32)),
 ('lembo', array([[1.0000004]], dtype=float32)),
 ('35-billion', array([[1.0000004]], dtype=float32)),
 ('thf', array([[1.0000004]], dtype=float32)),
 ('najd', array([[1.0000004]], dtype=float32)),
 ('meilan', array([[1.0000004]], dtype=float32)),
 ('silverleaf', array([[1.0000004]], dtype=float32)),
 ('featherbed', array([[1.0000004]], dtype=float32)),
 ('18-26', array([[1.0000004]], dtype=float32))]
```

```
similar_words(defining_model,'market',list(defining_model.keys()),20)
```

```
[('difference', array([[1.0000004]], dtype=float32)),
 ('forgotten', array([[1.0000004]], dtype=float32)),
 ('chester', array([[1.0000004]], dtype=float32)),
 ('protagonist', array([[1.0000004]], dtype=float32)),
 ('amounting', array([[1.0000004]], dtype=float32)),
 ('laude', array([[1.0000004]], dtype=float32)),
 ('congressionally', array([[1.0000004]], dtype=float32)),
 ('twinkling', array([[1.0000004]], dtype=float32)),
 ('dianna', array([[1.0000004]], dtype=float32)),
 ('replaceable', array([[1.0000004]], dtype=float32)),
 ('cnn.com', array([[1.0000004]], dtype=float32)),
 ('acupuncturist', array([[1.0000004]], dtype=float32)),
 ('lembo', array([[1.0000004]], dtype=float32)),
 ('35-billion', array([[1.0000004]], dtype=float32)),
 ('thf', array([[1.0000004]], dtype=float32)),
 ('najd', array([[1.0000004]], dtype=float32)),
 ('meilan', array([[1.0000004]], dtype=float32)),
```

```
        ('silverleaf', array([[1.0000004]], dtype=float32)),
        ('featherbed', array([[1.0000004]], dtype=float32)),
        ('18-26', array([[1.0000004]], dtype=float32))]


similar_words(defining_model,'stanford',list(defining_model.keys()),20)

     [('difference', array([[1.0000004]], dtype=float32)),
      ('forgotten', array([[1.0000004]], dtype=float32)),
      ('chester', array([[1.0000004]], dtype=float32)),
      ('protagonist', array([[1.0000004]], dtype=float32)),
      ('amounting', array([[1.0000004]], dtype=float32)),
      ('laude', array([[1.0000004]], dtype=float32)),
      ('congressionally', array([[1.0000004]], dtype=float32)),
      ('twinkling', array([[1.0000004]], dtype=float32)),
      ('dianna', array([[1.0000004]], dtype=float32)),
      ('replaceable', array([[1.0000004]], dtype=float32)),
      ('cnn.com', array([[1.0000004]], dtype=float32)),
      ('acupuncturist', array([[1.0000004]], dtype=float32)),
      ('lembo', array([[1.0000004]], dtype=float32)),
      ('35-billion', array([[1.0000004]], dtype=float32)),
      ('thf', array([[1.0000004]], dtype=float32)),
      ('najd', array([[1.0000004]], dtype=float32)),
      ('meilan', array([[1.0000004]], dtype=float32)),
      ('silverleaf', array([[1.0000004]], dtype=float32)),
      ('featherbed', array([[1.0000004]], dtype=float32)),
      ('18-26', array([[1.0000004]], dtype=float32))]



from nltk.tokenize import sent_tokenize, word_tokenize
import warnings

warnings.filterwarnings(action = 'ignore')

import gensim
from gensim.models import Word2Vec

def tsne_plot(model, word):
    "Creates and TSNE model and plots it"
    labels = []
    tokens = []
    if(word in defining_model.keys()):
```

```python
        for word in model.wv.vocab:
            tokens.append(model[word])
            labels.append(word)
        tsne_model = TSNE(perplexity=40, n_components=2, init='pca', n_iter=2500, random_state=23)
        new_values = tsne_model.fit_transform(tokens)

        x = []
        y = []
        for value in new_values:
            x.append(value[0])
            y.append(value[1])

        plt.figure(figsize=(16, 16))
        for i in range(len(x)):
            plt.scatter(x[i],y[i])
            plt.annotate(labels[i],
                        xy=(x[i], y[i]),
                        xytext=(5, 2),
                        textcoords='offset points',
                        ha='right',
                        va='bottom')
        plt.show()


f = open('/content/sample_data/vectors.txt', encoding="utf8")
model = gensim.models.Word2Vec(f)
for line in f:
    values = line.split()
    word = values[0]
    coefficient = np.asarray(values[1:], dtype='float32')
    model[word] = coefficient
f.close()


tsne_plot(model, 'life')
```

tsne_plot(model, 'market')

tsne_plot(model, 'stanford')