

```
1 ! pip install gensim
```

```
Requirement already satisfied: gensim in /usr/local/lib/python3.12/dist-packages (4.3.3)
Requirement already satisfied: numpy<2.0,>=1.18.5 in /usr/local/lib/python3.12/dist-packages (from gensim) (1.26.4)
Requirement already satisfied: scipy<1.14.0,>=1.7.0 in /usr/local/lib/python3.12/dist-packages (from gensim) (1.13.1)
Requirement already satisfied: smart-open>=1.8.1 in /usr/local/lib/python3.12/dist-packages (from gensim) (7.3.1)
Requirement already satisfied: wrapt in /usr/local/lib/python3.12/dist-packages (from smart-open>=1.8.1->gensim) (1.17.3)
```

Task 1

```
1 import numpy as np
2 import gensim.downloader as api
3
4
5
6 class WordEmbedder:
7     def __init__(self, model_name: str = 'glove-wiki-gigaword-50'):
8         print(f"Loading model '{model_name}'")
9         self.model = api.load(model_name)
10        self.vector_size = self.model.vector_size
11        print(f"Model '{model_name}' loaded successfully! Vector size: {self.vector_size}")
12
13    def get_vector(self, word: str):
14        if word in self.model.key_to_index:
15            return self.model[word]
16        else:
17            print(f"Warning: '{word}' not in vocabulary (OOV).")
18            return None
19
20    def get_similarity(self, word1: str, word2: str):
21        if word1 in self.model.key_to_index and word2 in self.model.key_to_index:
22            return self.model.similarity(word1, word2)
23        else:
24            print(f"Warning: One or both words are OOV: '{word1}', '{word2}'.")
25            return None
26
27    def get_most_similar(self, word: str, top_n: int = 10):
28        if word in self.model.key_to_index:
29            return self.model.most_similar(word, topn=top_n)
30        else:
31            print(f"Warning: '{word}' not in vocabulary (OOV).")
32            return []
33
34    def embed_document(self, document: str):
35        from nltk.tokenize import word_tokenize
36        import nltk
37        nltk.download('punkt_tab', quiet=True)
38
39        tokens = word_tokenize(document.lower())
40        vectors = []
41
42        for token in tokens:
43            if token in self.model.key_to_index:
44                vectors.append(self.model[token])
45
46        if not vectors:
47            print("No known words in document, returning zero vector.")
48            return np.zeros(self.vector_size)
49
50        return np.mean(vectors, axis=0)
51
```

```
1 def main():
2     embedder = WordEmbedder('glove-wiki-gigaword-50')
3
4     print("\nVector for 'king':")
5     print(embedder.get_vector('king'))
6
7     print("\nSimilarity between 'king' and 'queen':", embedder.get_similarity('king', 'queen'))
8     print("Similarity between 'king' and 'man':", embedder.get_similarity('king', 'man'))
9
10    print("\nMost similar to 'computer':")
11    similar_words = embedder.get_most_similar('computer')
12    for word, score in similar_words:
13        print(f"{word}: {score:.4f}")
14    sentence = "The queen rules the country."
```

```

15     print(f"\nEmbedding for sentence: \'{sentence}\'")
16     doc_vector = embedder.embed_document(sentence)
17     print(doc_vector)
18     print("Vector shape:", doc_vector.shape)
19
20
21 if __name__ == "__main__":
22     main()
23

```

Loading model 'glove-wiki-gigaword-50'
[=====] 100.0% 66.0/66.0MB downloaded
Model 'glove-wiki-gigaword-50' loaded successfully! Vector size: 50

Vector for 'king':
[0.50451 0.68607 -0.59517 -0.022801 0.60046 -0.13498 -0.08813
 0.47377 -0.61798 -0.31012 -0.076666 1.493 -0.034189 -0.98173
 0.68229 0.81722 -0.51874 -0.31503 -0.55809 0.66421 0.1961
 -0.13495 -0.11476 -0.30344 0.41177 -2.223 -1.0756 -1.0783
 -0.34354 0.33505 1.9927 -0.04234 -0.64319 0.71125 0.49159
 0.16754 0.34344 -0.25663 -0.8523 0.1661 0.40102 1.1685
 -1.0137 -0.21585 -0.15155 0.78321 -0.91241 -1.6106 -0.64426
 -0.51042]

Similarity between 'king' and 'queen': 0.78390425
Similarity between 'king' and 'man': 0.53093773

Most similar to 'computer':
computers: 0.9165
software: 0.8815
technology: 0.8526
electronic: 0.8126
internet: 0.8060
computing: 0.8026
devices: 0.8016
digital: 0.7992
applications: 0.7913
pc: 0.7883

Embedding for sentence: "The queen rules the country."
[0.04564168 0.36530998 -0.55974334 0.04014383 0.09655549 0.15623933
 -0.33622834 -0.12495166 -0.01031508 -0.5006717 0.18690467 0.17482166
 -0.268985 -0.03096624 0.36686516 0.29983264 0.01397333 -0.06872118
 -0.3260683 -0.210115 0.16835399 -0.03151734 -0.06204716 0.04301083
 -0.06958768 -1.7792168 -0.54365396 -0.06104483 -0.17618 0.009181
 3.3916333 0.08742473 -0.4675417 -0.213435 0.02391887 -0.04470453
 0.20636833 -0.12902866 -0.28527132 -0.2431805 -0.3114423 -0.03833717
 0.11977985 -0.01418401 -0.37086335 0.22069354 -0.28848937 -0.36188802
 -0.00549529 -0.46997246]
Vector shape: (50,)

Task 2 Document embedding

```

1 sentence = 'The quick brown fox jumps over the lazy dog.'
2 embedder = WordEmbedder()
3 doc_vector = embedder.embed_document(sentence)
4 print(doc_vector)
5 print("Vector shape:", doc_vector.shape)


```

Loading model 'glove-wiki-gigaword-50'
Model 'glove-wiki-gigaword-50' loaded successfully! Vector size: 50
[0.018704 -0.01994951 -0.152806 -0.05444079 0.267892 0.27261212
 -0.594531 -0.1579832 0.05074116 -0.22430201 -0.19797495 0.07773571
 -0.18965401 0.23584989 0.2018865 -0.03365446 0.1494343 -0.123215
 -0.77456605 -0.400587 -0.1006896 0.15792981 0.37431902 -0.02621669
 0.20702538 -1.328874 -0.2806339 0.34456852 0.22626872 -0.285655
 2.6676202 0.18003204 -0.23781231 0.09433399 -0.04113727 -0.09008334
 -0.0373017 -0.0350816 0.08468361 -0.4464623 -0.15689898 0.25797132
 -0.33406436 0.11145066 0.040368 -0.18526098 0.17259666 -0.26326978
 0.18413429 0.19813938]
Vector shape: (50,)

Task 3

```

1 from gensim.models import Word2Vec
2 from gensim.utils import simple_preprocess
3 from pathlib import Path
4
5
6 class Word2VecTrainer:
7     def __init__(self, data_path: str, model_path: str = "models/word2vec.model"):
```

```

8     self.data_path = Path(data_path)
9     self.model_path = Path(model_path)
10    self.model = None
11
12    def load_data(self):
13        print(f"Loading data from: {self.data_path}")
14        sentences = []
15        with open(self.data_path, "r", encoding="utf-8") as f:
16            for line in f:
17                tokens = simple_preprocess(line)
18                if tokens:
19                    sentences.append(tokens)
20        print(f"Loaded {len(sentences)} sentences.")
21        return sentences
22
23    def train(self, vector_size=50, window=5, min_count=1, workers=4, epochs=10):
24        sentences = self.load_data()
25        print("Training Word2Vec model...")
26        self.model = Word2Vec(
27            sentences=sentences,
28            vector_size=vector_size,
29            window=window,
30            min_count=min_count,
31            workers=workers,
32            epochs=epochs,
33        )
34        print("Training completed.")
35
36    def save_model(self):
37        self.model_path.parent.mkdir(parents=True, exist_ok=True)
38        self.model.save(str(self.model_path))
39        print(f"Model saved to {self.model_path}")
40
41    def load_model(self):
42        print(f"Loading model from {self.model_path}")
43        self.model = Word2Vec.load(str(self.model_path))
44        print("Model loaded successfully.")
45        return self.model
46
47    def get_similarity(self, w1, w2):
48        if w1 in self.model.wv and w2 in self.model.wv:
49            return self.model.wv.similarity(w1, w2)
50        else:
51            print(f"One of the words not in vocabulary: {w1}, {w2}")
52            return None
53
54    def get_most_similar(self, word, top_n=5):
55        if word in self.model.wv:
56            return self.model.wv.most_similar(word, topn=top_n)
57        else:
58            print(f"Word not in vocabulary: {word}")
59            return []
60

```

```

1 trainer = Word2VecTrainer(data_path="/content/en_ewt-ud-train.txt")
2 trainer.train(vector_size=50, window=5, min_count=1, epochs=20)
3
4 trainer.save_model()
5
6 trainer.load_model()
7 W1 = 'king'
8 W2 = 'queen'
9 print(f"\nTừ tương tự '{W1}':")
10 print(trainer.get_most_similar(W1))
11
12 print(f"\nĐộ tương đồng giữa '{W1}' và '{W2}':")
13 print(trainer.get_similarity(W1, W2))

```

Loading data from: /content/en_ewt-ud-train.txt
 Loaded 14196 sentences.
 Training Word2Vec model...
 Training completed.
 Model saved to models/word2vec.model
 Loading model from models/word2vec.model
 Model loaded successfully.

Từ tương tự 'king':
 [('barely', 0.9760759472846985), ('goldman', 0.9746447801589966), ('welcome', 0.9742516875267029), ('bowtie', 0.9677373766)

```
Độ tương đồng giữa 'king' và 'queen':
0.8856816
```

Task 4

```

1 import os
2 import re
3 from pyspark.sql import SparkSession
4 from pyspark.sql.functions import col, lower, regexp_replace, split
5 from pyspark.ml.feature import Word2Vec
6
7 os.environ["HADOOP_HOME"] = "C:\\\\hadoop"
8 os.environ["SPARK_LOCAL_DIRS"] = "C:\\\\spark-temp"
9 os.makedirs("C:\\\\hadoop\\\\bin", exist_ok=True)
10 os.makedirs("C:\\\\spark-temp", exist_ok=True)
11
12 spark = (
13     SparkSession.builder.appName("Lab4_Spark_Word2Vec_C4")
14     .config("spark.driver.memory", "4g")
15     .config("spark.executor.memory", "4g")
16     .getOrCreate()
17 )
18
19 data_path = "/content/c4-train.00000-of-01024-30K.json"
20
21 try:
22     df = spark.read.json(data_path)
23     df_clean = (
24         df.select("text")
25         .withColumn("text", lower(col("text")))
26         .withColumn("text", regexp_replace(col("text"), "[^a-zA-Z\\s]", ""))
27         .withColumn("words", split(col("text"), "\\s+"))
28         .filter(col("words").isNotNull())
29     )
30
31     word2Vec = Word2Vec(
32         vectorSize=100,
33         minCount=5,
34         inputCol="words",
35         outputCol="result"
36     )
37
38     model = word2Vec.fit(df_clean)
39
40     word = "computer"
41     synonyms = model.findSynonyms(word, 5)
42     print(f"Các từ giống '{word}' nhất (từ model Spark):")
43     synonyms.show()
44
45 except Exception as e:
46     print(f"Đã có lỗi xảy ra: {e}")
47
48 finally:
49     spark.stop()
50

```

```
Các từ giống 'computer' nhất (từ model Spark):
+-----+-----+
|   word|      similarity|
+-----+-----+
|stationery|0.6567521691322327|
|  update|0.6258612275123596|
| desktop| 0.621239423751831|
| browser|0.6166326403617859|
|    laptop|0.6047263741493225|
+-----+-----+
```

Task 5

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 from sklearn.decomposition import PCA
4 from sklearn.manifold import TSNE
5 import numpy as np
6
7 words = [

```

```

8     'engineer', 'artist',
9     'newyork', 'london',
10    'facebook', 'tesla',
11    'bird', 'horse',
12    'love', 'fear',
13    'phone', 'computer', 'movie'
14 ]
15
16
17 word_vectors = np.array([embedder.get_vector(w) for w in words])
18 pca = PCA(n_components=2)
19 vectors_pca = pca.fit_transform(word_vectors)
20
21
22 tsne = TSNE(n_components=2, random_state=42, perplexity=5)
23 vectors_tsne = tsne.fit_transform(word_vectors)
24
25 plt.style.use('seaborn-v0_8-whitegrid')
26 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(20, 9))
27
28 ax1.scatter(vectors_pca[:, 0], vectors_pca[:, 1], c='blue', edgecolors='k', alpha=0.6)
29 for i, word in enumerate(words):
30     ax1.annotate(word, (vectors_pca[i, 0], vectors_pca[i, 1]), fontsize=12)
31 ax1.set_title("Trục quan hóa bằng PCA", fontsize=16)
32 ax1.set_xlabel("Thành phần chính 1", fontsize=12)
33 ax1.set_ylabel("Thành phần chính 2", fontsize=12)
34 ax1.grid(True)
35
36 ax2.scatter(vectors_tsne[:, 0], vectors_tsne[:, 1], c='green', edgecolors='k', alpha=0.6)
37 for i, word in enumerate(words):
38     ax2.annotate(word, (vectors_tsne[i, 0], vectors_tsne[i, 1]), fontsize=12)
39 ax2.set_title("Trục quan hóa bằng t-SNE", fontsize=16)
40 ax2.set_xlabel("t-SNE dimension 1", fontsize=12)
41 ax2.set_ylabel("t-SNE dimension 2", fontsize=12)
42 ax2.grid(True)
43
44 plt.show()
45

```

```

/usr/local/lib/python3.12/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 7847 (\N{LATIN SMALL LETTER A W
fig.canvas.print_figure(bytes_io, **kw)
/usr/local/lib/python3.12/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 7921 (\N{LATIN SMALL LETTER U W
fig.canvas.print_figure(bytes_io, **kw)
/usr/local/lib/python3.12/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 7857 (\N{LATIN SMALL LETTER A W
fig.canvas.print_figure(bytes_io, **kw)

```



