

# Experiment report: TextAnalysis

Done by student: **Bylkova Kristina (伯汀娜), 1820249049**

Course: **Big Data Analysis Technology**

Date: **October 2024**

## Task

1. Learn and understand the word embedding, draw a diagram to show the mechanism and work flow of word embedding;
2. Finish the experiment and write the experiment report.

## Word embedding

Word embedding is a technique which converts a text into a numerical representation, mapping each word to a vector. This vector captures semantic meaning and distinctive features of every word. The main feature is that the closer words are in the vector space, more similar they are expected to be in meaning. This method is widely used in text analysis.

The work principle is shown in the diagram below.

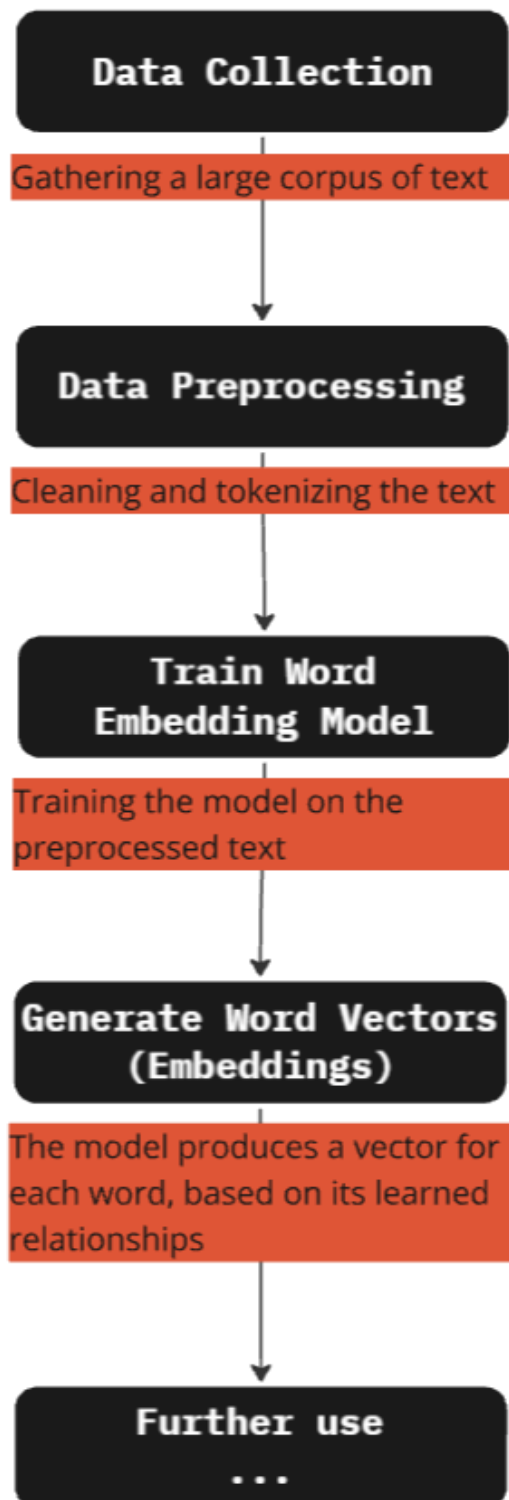


Figure 1: Diagram of word embedding

## Code

```
1 import tensorflow as tf
2 from sklearn.metrics import accuracy_score
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.model_selection import train_test_split
5 from sklearn.feature_extraction.text import TfidfVectorizer
6 import matplotlib.pyplot as plt
7 import time
8 import numpy as np
9 import pandas as pd
10 import nltk
11 from nltk.corpus import movie_reviews
12 import tensorflow as tf
13
14 # Tokenizer from Keras, used to convert text into sequences
15 tokenizer = tf.keras.preprocessing.text.Tokenizer()
16
17 # Using Keras API components from TensorFlow for model creation
18 pad_sequences = tf.keras.preprocessing.sequence.pad_sequences
19 Embedding = tf.keras.layers.Embedding
20 LSTM = tf.keras.layers.LSTM
21 Dense = tf.keras.layers.Dense
22 Sequential = tf.keras.models.Sequential
23
24 # Downloading movie reviews dataset from NLTK
25 # nltk.download('movie_reviews')
26
27 # Converting the raw text data into a structured DataFrame with two columns: 'text' and 'label'
28 # 'text' contains the review, and 'label' is the sentiment (positive/negative)
29 documents = [(list(movie_reviews.words(fileid)), category)
30              for category in movie_reviews.categories()
31              for fileid in movie_reviews.fileids(category)]
32 df = pd.DataFrame(documents, columns=['text', 'label'])
33 # Mapping labels to binary values (1 for positive, 0 for negative)
34 df['label'] = df['label'].map({'pos': 1, 'neg': 0})
35 # Joining the review words into a single string
36 df['text'] = df['text'].apply(lambda x: ' '.join(x))
37
38 # TF-IDF vectorization: transforming the text data into numerical features
39 vectorizer = TfidfVectorizer(max_features=5000) # Using max 5000 features
40 # Transforming text into TF-IDF feature vectors
41 X_tfidf = vectorizer.fit_transform(df['text']).toarray()
42 y = df['label'].values # Extracting labels
43
44 # Splitting the data into training and testing sets
45 X_train_tfidf, X_test_tfidf, y_train, y_test = train_test_split(
46     X_tfidf, y, test_size=0.2, random_state=42)
47
48 # Training a Logistic Regression model using the TF-IDF features
49 start_time = time.time()
50 model_tfidf = LogisticRegression() # Initializing the model
51 model_tfidf.fit(X_train_tfidf, y_train) # Training the model
```

```

52 training_time_tfidf = time.time() - start_time # Calculating training time
53 # Making predictions on the test set
54 y_pred_tfidf = model_tfidf.predict(X_test_tfidf)
55 accuracy_tfidf = accuracy_score(
56     y_test, y_pred_tfidf) # Calculating the accuracy
57 print(f"Accuracy of TF-IDF model: {accuracy_tfidf*100:.2f}%")
58
59 # Tokenizing the text for LSTM
60 tokenizer = tf.keras.preprocessing.text.Tokenizer(
61     num_words=5000) # Using max 5000 words
62 tokenizer.fit_on_texts(df['text']) # Fitting the tokenizer on the dataset
63 # Converting text to sequences of integers
64 X_seq = tokenizer.texts_to_sequences(df['text'])
65 X_pad = tf.keras.preprocessing.sequence.pad_sequences(
66     X_seq, maxlen=500) # Padding sequences to ensure equal length
67
68 # Loading GloVe embeddings for the LSTM model
69 embedding_index = {}
70 with open('glove.6B.100d.txt', 'r', encoding='utf-8') as f:
71     for line in f:
72         values = line.split()
73         word = values[0] # The word
74         coefs = np.asarray(values[1:], dtype='float32') # The embedding vector
75         embedding_index[word] = coefs # Storing the word and its vector
76
77 # Creating an embedding matrix where each word is represented by its GloVe embedding
78 word_index = tokenizer.word_index
79 # Matrix size of (vocab_size, embedding_dim)
80 embedding_matrix = np.zeros((5000, 100))
81 for word, i in word_index.items():
82     if i < 5000:
83         embedding_vector = embedding_index.get(word)
84         if embedding_vector is not None:
85             # Filling in the embedding matrix
86             embedding_matrix[i] = embedding_vector
87
88 # Building the LSTM model
89 model_lstm = Sequential() # Initializing a sequential model
90 model_lstm.add(Embedding(input_dim=5000, output_dim=100, weights=[
91     embedding_matrix], trainable=False)) # Embedding layer with GloVe weights
92 # LSTM layer with dropout for regularization
93 model_lstm.add(LSTM(100, dropout=0.2, recurrent_dropout=0.2))
94 # Dense output layer with sigmoid activation for binary classification
95 model_lstm.add(Dense(1, activation='sigmoid'))
96
97 # Compiling the LSTM model
98 model_lstm.compile(loss='binary_crossentropy',
99     optimizer='adam', metrics=['accuracy'])
100
101 # Splitting the padded sequences into training and testing sets
102 X_train_pad, X_test_pad, y_train, y_test = train_test_split(
103     X_pad, y, test_size=0.2, random_state=42)

```

```

104
105 # Training the LSTM model
106 start_time = time.time()
107 model_lstm.fit(X_train_pad, y_train, validation_data=(
108     X_test_pad, y_test), epochs=15, batch_size=128)
109     # X_test_pad, y_test), epochs=3, batch_size=128)
110 # Calculating training time for LSTM
111 training_time_lstm = time.time() - start_time
112 # Evaluating the LSTM model
113 scores = model_lstm.evaluate(X_test_pad, y_test, verbose=0)
114 accuracy_lstm = scores[1] # Extracting accuracy
115 print(f"Accuracy of LSTM: {accuracy_lstm*100:.2f}%")
116
117 # Results comparison between TF-IDF + Logistic Regression and LSTM + GloVe
118 results = pd.DataFrame({
119     'Model': ['TF-IDF + Logistic Regression', 'LSTM + GloVe'],
120     'Accuracy': [accuracy_tfidf, accuracy_lstm],
121     'Training Time (seconds)': [training_time_tfidf, training_time_lstm]
122 })
123
124 print(results)
125
126 # Plotting accuracy comparison
127 plt.figure(figsize=(10, 5))
128 plt.bar(results['Model'], results['Accuracy'], color=['blue', 'green'])
129 plt.title('Accuracy Comparison')
130 plt.ylabel('Accuracy')
131 plt.show()
132
133 # Plotting training time comparison
134 plt.figure(figsize=(10, 5))
135 plt.bar(results['Model'], results['Training Time (seconds)'],
136     color=['blue', 'green'])
137 plt.title('Training Time Comparison')
138 plt.ylabel('Training Time (seconds)')
139 plt.show()

```

## Experiment

In this experiment after the code execution we get the accuracy and training time output for each model (TF-IDF+ Logistic regression and LSTM + GloVe). The accuracy of TF-IDF: 80.75%.

At first, I ran the code with the number of epochs = 15. The accuracy of LSTM for this case is 75.00%. It is shown that the accuracy is higher using TF-IDF + Logistic Regression than LSTM + GloVe. But the training time is much higher for LSTM + GloVe.

```
PS C:\Users\Admin\Desktop\TextAnalysis> python -u "c:\Users\Admin\Desktop\TextAnalysis\experiment.py"
Accuracy of TF-IDF model: 80.75%
2024-10-24 14:12:59.104465: I tensorflow/core/platform/cpu_feature_guard.cc:182] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.
To enable the following instructions: SSE SSE2 SSE3 SSE4.1 SSE4.2 AVX AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
Epoch 1/15
13/13 [=====] - 68s 5s/step - loss: 0.6979 - accuracy: 0.5088 - val_loss: 0.6870 - val_accuracy: 0.5450
Epoch 2/15
13/13 [=====] - 69s 5s/step - loss: 0.6671 - accuracy: 0.6069 - val_loss: 0.6803 - val_accuracy: 0.5600
Epoch 3/15
13/13 [=====] - 68s 5s/step - loss: 0.6471 - accuracy: 0.6338 - val_loss: 0.6485 - val_accuracy: 0.6400
Epoch 4/15
13/13 [=====] - 91s 7s/step - loss: 0.6233 - accuracy: 0.6481 - val_loss: 0.6446 - val_accuracy: 0.6425
- accuracy: 0.7269 - val_loss: 0.5539 - val_accuracy: 0.7300
Epoch 9/15
13/13 [=====] - 74s 6s/step - loss: 0.5436 - accuracy: 0.7294 - val_loss: 0.6199 - val_accuracy: 0.6750
Epoch 10/15
13/13 [=====] - 73s 6s/step - loss: 0.5065 - accuracy: 0.7675 - val_loss: 0.5601 - val_accuracy: 0.7550
Epoch 11/15
13/13 [=====] - 75s 6s/step - loss: 0.5012 - accuracy: 0.7600 - val_loss: 0.5757 - val_accuracy: 0.7375
Epoch 12/15
13/13 [=====] - 75s 6s/step - loss: 0.4964 - accuracy: 0.7606 - val_loss: 0.5557 - val_accuracy: 0.7350
Epoch 13/15
13/13 [=====] - 77s 6s/step - loss: 0.4756 - accuracy: 0.7800 - val_loss: 0.6296 - val_accuracy: 0.6950
Epoch 14/15
13/13 [=====] - 76s 6s/step - loss: 0.4934 - accuracy: 0.7638 - val_loss: 0.6571 - val_accuracy: 0.6400
Epoch 15/15
13/13 [=====] - 88s 7s/step - loss: 0.4870 - accuracy: 0.7631 - val_loss: 0.5449 - val_accuracy: 0.7500
Accuracy of LSTM: 75.00%
Model Accuracy Training Time (seconds)
0 TF-IDF + Logistic Regression 0.8075 0.170569
1 LSTM + GloVe 0.7500 1122.635486
```

Figure 2: Code execution: 15 epochs

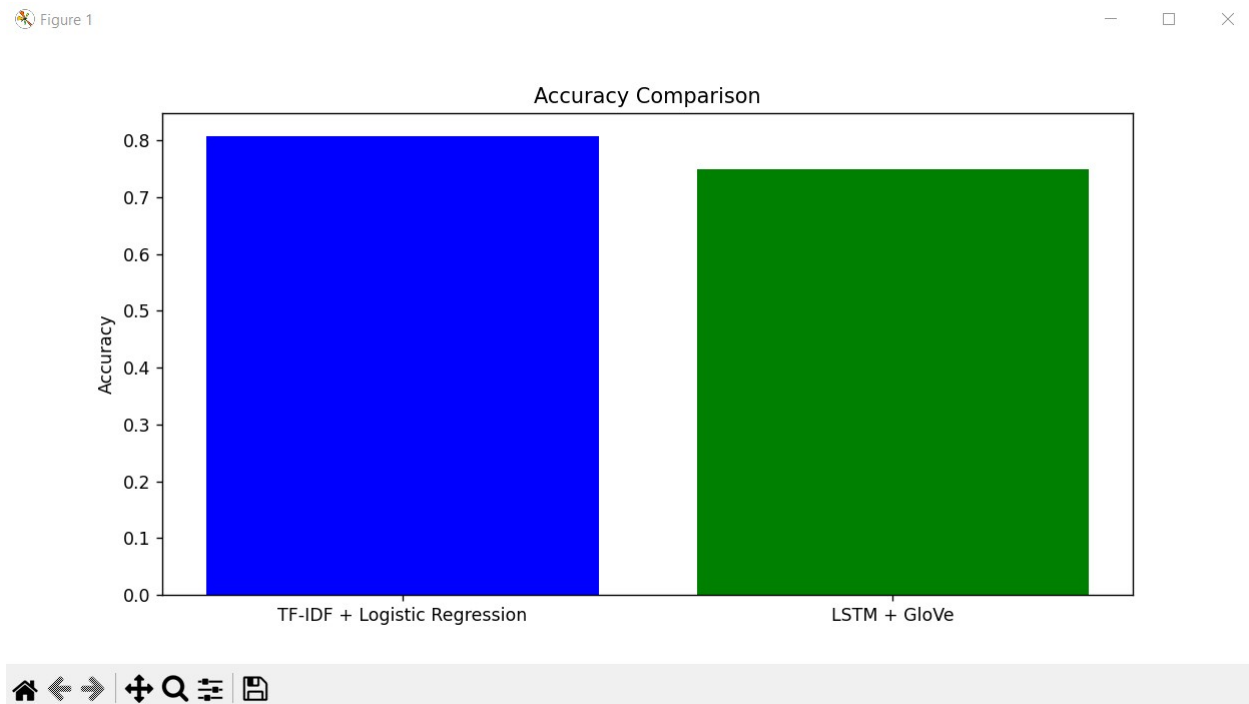


Figure 3: Accuracy comparison: 15 epochs



Figure 4: Training time comparison: 15 epochs

Then I ran the code with the number of epochs = 3. The accuracy of LSTM for this case is 61.00%. The situations with the accuracy and the training time is the same.

```
PS C:\Users\Admin\Desktop\TextAnalysis> python -u "c:\Users\Admin\Desktop\TextAnalysis\experiment.py"
Accuracy of TF-IDF model: 80.75%
2024-10-24 14:34:53.687950: I tensorflow/core/platform/cpu_feature_guard.cc:182] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.
To enable the following instructions: SSE SSE2 SSE3 SSE4.1 SSE4.2 AVX AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
Epoch 1/3
13/13 [=====] - 58s 4s/step - loss: 0.6869 - accuracy: 0.5394 - val_loss: 0.6822 - val_accuracy: 0.5375
Epoch 2/3
13/13 [=====] - 56s 4s/step - loss: 0.6630 - accuracy: 0.6006 - val_loss: 0.6674 - val_accuracy: 0.5800
Epoch 3/3
13/13 [=====] - 53s 4s/step - loss: 0.6478 - accuracy: 0.6244 - val_loss: 0.6517 - val_accuracy: 0.6100
Accuracy of LSTM: 61.00%


|   | Model                        | Accuracy | Training Time (seconds) |
|---|------------------------------|----------|-------------------------|
| 0 | TF-IDF + Logistic Regression | 0.8075   | 0.174005                |
| 1 | LSTM + GloVe                 | 0.6100   | 168.287405              |


PS C:\Users\Admin\Desktop\TextAnalysis>
```

Figure 5: Code execution: 3 epochs

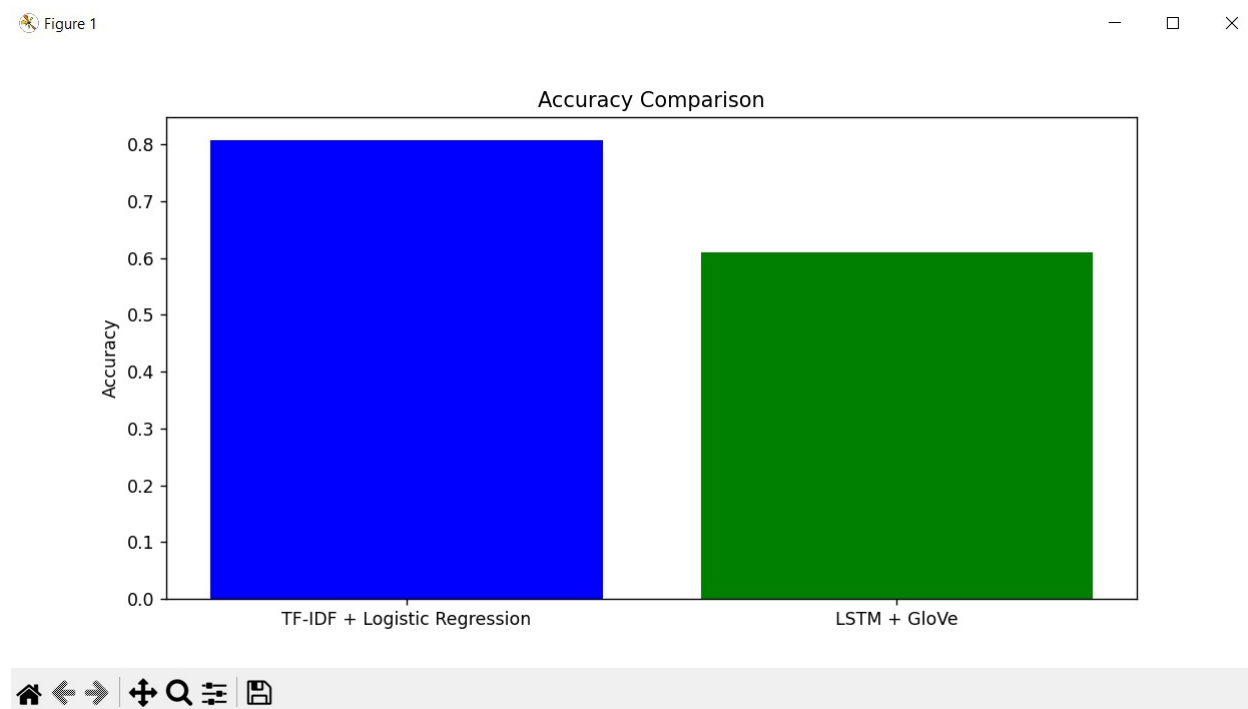


Figure 6: Accuracy comparison: 3 epochs





Figure 7: Training time comparison: 3 epochs