

جهت سهولت در نوشتن کد ها، هر فاز به صورت جداگانه نوشته شده و در نهایت خروجی فاز برای استفاده در مراحل بعدی ذخیره شده.

کتابخانه های استفاده شده:

```
1. import pandas as pd
2. from sklearn.model_selection import train_test_split
3. from textblob import TextBlob
4. import matplotlib.pyplot as plt
5. from nltk.stem import PorterStemmer
6. import re
7.
8. from nltk.corpus import stopwords
9. from sklearn.feature_extraction.text import CountVectorizer
10. from sklearn.cluster import KMeans
11. from sklearn.manifold import TSNE
12. import matplotlib.pyplot as plt
13. import numpy as np
14.
15. from tensorflow.keras.preprocessing.text import Tokenizer
16. from tensorflow.keras.preprocessing.sequence import pad_sequences
17.
18. import urllib.request
19. import zipfile
20. import os
21.
22. from tensorflow.keras.models import Sequential
23. from tensorflow.keras.layers import Embedding, GRU, Dense, Bidirectional
24. from tensorflow.keras.callbacks import History
```

## ۱-پیش پردازش

در مرحله اول یک تابع برای بازکردن فایل اکسل نوشته شده:

```
1. # Step 1: Load the dataset
2. def load_dataset(file_path):
3.     df = pd.read_csv(file_path)
4.     return df
```

### مرحله ۱-۱-

در این مرحله مثبت ها را به یک تبدیل کردیم و منفی ها را به صفر

```
1. # Step 1.1: Convert positive labels to 1 and negative labels to 0
2. def convert_labels(df):
3.     df['sentiment'] = df['sentiment'].apply(lambda x: 1 if x == 'positive' else 0)
4.     return df
```

### مرحله ۲-۱-

در این مرحله رکود های تکراری حذف شده. از inplace استفاده شده تا فایل جدیدی ایجاد نشود و در همان فایل تغییرات اعمال شود.

```
1. # Step 1.2: Remove duplicate rows
2. def remove_duplicates(df):
3.     df.drop_duplicates(inplace=True)
```

مرحله ۳-۱-

در این مرحله ۳۰ درصد داده ها برای آموزش با رندم استیت ۴۲ جدا شده اند:

```
1. # Step 1.3: Randomly select 30% of the dataset
2. def random_selection(df):
3.     df = df.sample(frac=0.3, random_state=42)
4.     return df
```

مرحله ۴-۱-

در این مرحله یک ستون جدید ایجاد شده و با استفاده از کتابخانه textBlob قطبیت جمله در آن ثبت شده:

```
1. # Step 1.4: Calculate polarity using TextBlob
2. def get_polarity(df):
3.     print('step-1-4')
4.     df['polarity'] = df['review'].apply(lambda x: TextBlob(str(x)).sentiment.polarity)
5.     return df
```

مرحله ۵-۱-

در این مرحله از ستون sentiment ستون های ۱ را مثبت و ۰ را منفی در نظر گرفتیم و با پلات نمودار رسم کرده و ذخیره کردیم:

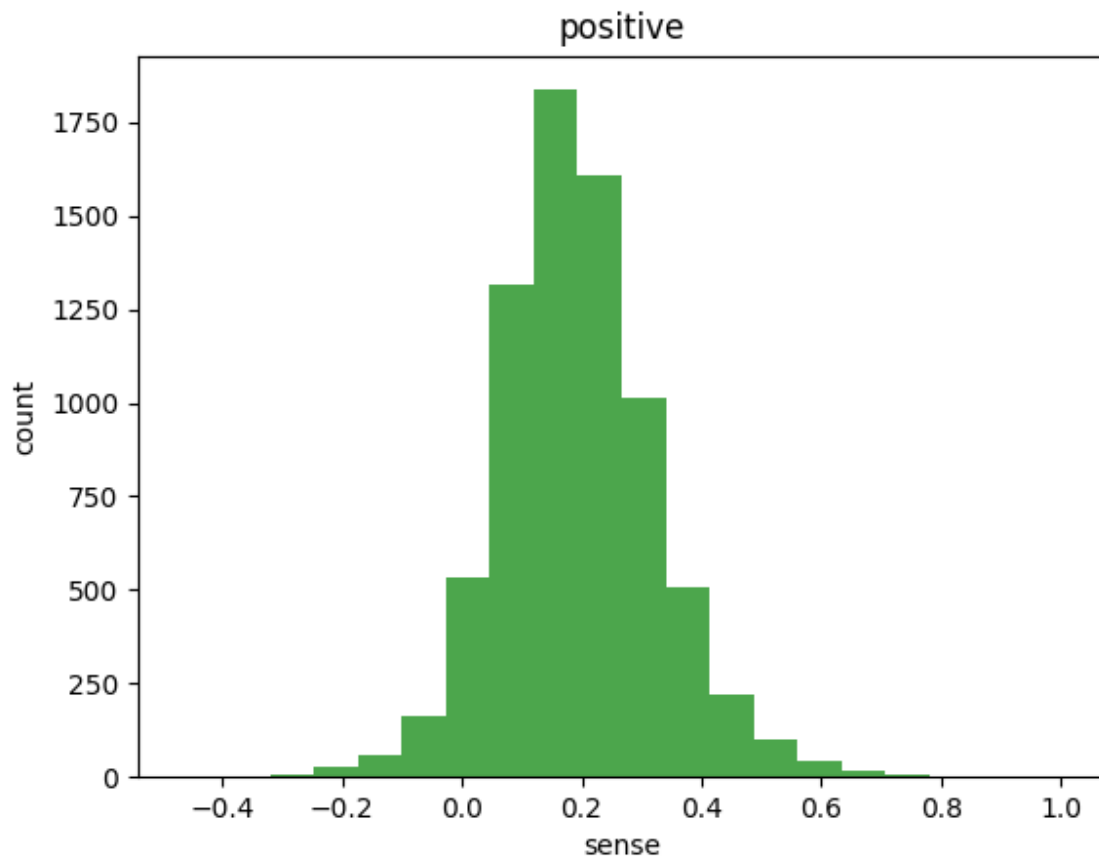
```
1. # Step 1.5: Separate positive and negative polarities, plot histograms, and analyze means
2. def plot_sentiment_histograms(df):
3.     print('step-1-5')
4.     positive_reviews = df[df['sentiment'] == 1]
5.     negative_reviews = df[df['sentiment'] == 0]
6.
7.     # Plotting histogram for positive sentiment
8.     plt.hist(positive_reviews['polarity'], bins=20, color='green', alpha=0.7)
9.     plt.title('positive')
10.    plt.xlabel('sense')
11.    plt.ylabel('count')
12.    plt.savefig('positive.png')
13.    plt.show()
14.
15.    # Plotting histogram for negative sentiment
16.    plt.hist(negative_reviews['polarity'], bins=20, color='red', alpha=0.7)
17.    plt.title('negative')
18.    plt.xlabel('sense')
```

```

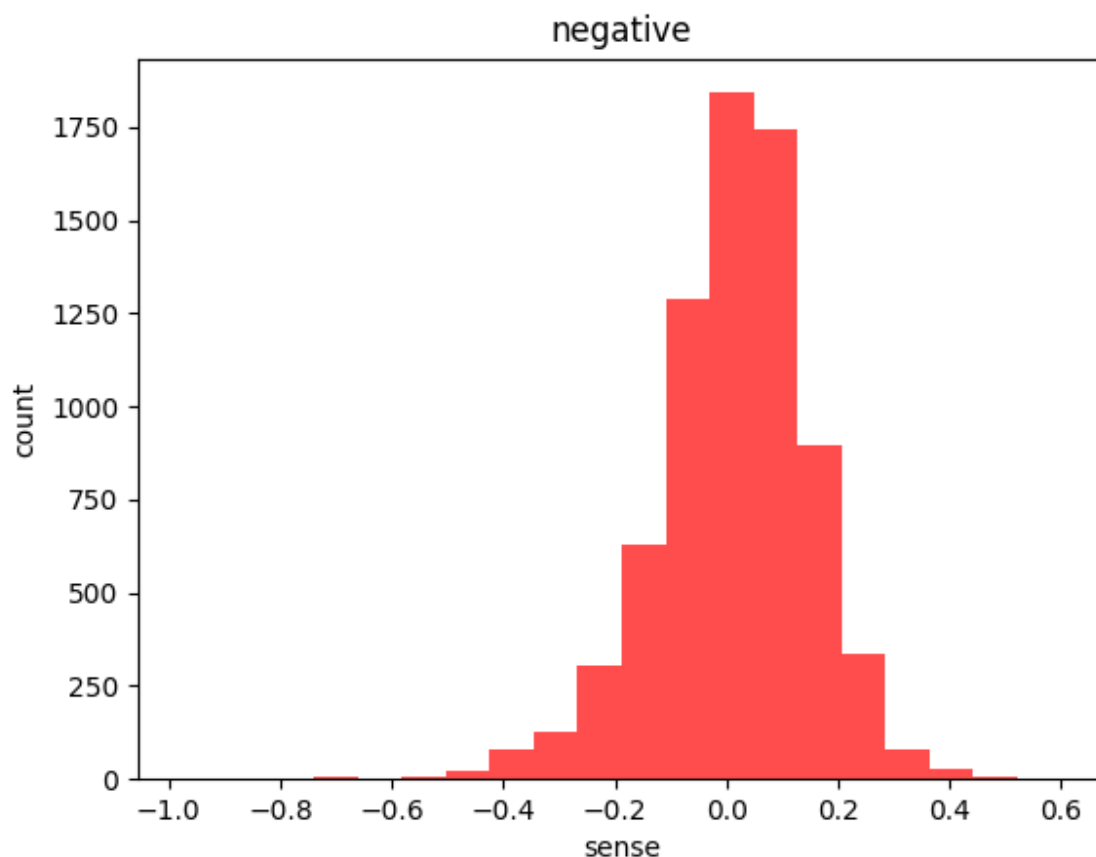
19. plt.ylabel('count')
20. plt.savefig('negative.png')
21. plt.show()
22.
23. # Calculate mean polarity for positive and negative reviews
24. mean_polarity_positive = positive_reviews['polarity'].mean()
25. mean_polarity_negative = negative_reviews['polarity'].mean()
26.
27. # Print mean polarities
28. print("Mean Polarity for Positive Reviews:", mean_polarity_positive)
29. print("Mean Polarity for Negative Reviews:", mean_polarity_negative)
30.

```

خروجی مثبت:



خروجی منفی:



خروجی میانگین ها در کنسول:

```
1. Mean Polarity for Positive Reviews: 0.19375026016013788
2. Mean Polarity for Negative Reviews: 0.012643524175639971
```

مرحله ۱-۶-

در این مرحله ابتدا تمام حروف کوچک شده و با رجکس اصلاح شده اند سپس استاپ ورد ها حذف شده و  
استم کلمات استخراج شده:

```
1. # Step 1.6: Apply preprocessing steps
2. def preprocess_text(text):
3.     text = text.lower()
4.     text = re.sub(r'\W', ' ', text)
5.     text = re.sub(r'\s+', ' ', text)
6.
7.     stop_words = set(stopwords.words('english'))
8.     text = ' '.join([word for word in text.split() if word not in stop_words])
9.
10.    stemmer = PorterStemmer()
```

```
11. text = ' '.join([stemmer.stem(word) for word in text.split()])
```

مرحله ۷-۱-

برای تقسیم داده به شکر زیر اقدام شده:

```
1. # Step 1.7: Split the dataset into training and testing sets
2. def split_dataset(df):
3.     train_data, test_data, train_labels, test_labels = train_test_split(
4.         df['review'], df['sentiment'], test_size=0.5, random_state=42
5.     )
6.     return train_data, test_data, train_labels, test_labels
```

۲-خوشه بندی

مرحله ۲-۱-

در این CountVectorizer فراخوانی شده و داده های آموزشی و تست فیت ترنسفورم شده اند.

```
1. # Step 2-1: Create Bag of Words (BOW) representation for each review
2. def create_bow_representation(reviews_train, reviews_test):
3.     vectorizer = CountVectorizer()
4.
5.     bow_matrix_train = vectorizer.fit_transform(reviews_train)
6.     bow_matrix_test = vectorizer.transform(reviews_test)
7.
8.     return bow_matrix_train, bow_matrix_test
```

مرحله ۲-۲-

در این مرحله خوشه بند k-means را فراخوانی کرده و دسته بندی را انجام دادیم:

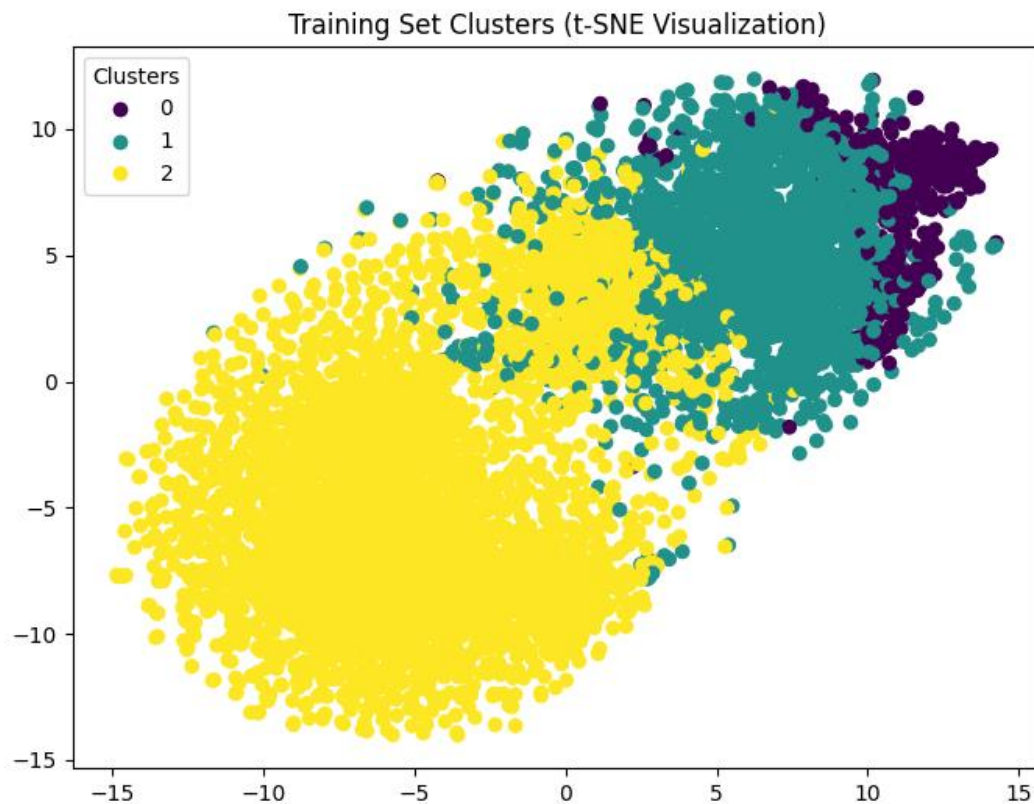
```
1. # Step 2-2: Cluster the training and testing sets using K-means
2. def kmeans_clustering(bow_matrix_train, bow_matrix_test, num_clusters=3):
3.     kmeans = KMeans(n_clusters=num_clusters, random_state=42)
4.
5.     # Fit on training set
6.     train_clusters = kmeans.fit_predict(bow_matrix_train)
7.
8.     # Predict on testing set
9.     test_clusters = kmeans.predict(bow_matrix_test)
10.
11.     return train_clusters, test_clusters
12.
```

مرحله ۲-۳-

در این مرحله با استفاده tsne خوشه ها رسم شده اند و در پلات نمایش داده شده اند:

```
1. # Step 2-3: Visualize clusters using t-sne
2. def visualize_clusters(bow_matrix, clusters, title):
3.     print('step-2-3')
4.     tsne = TSNE(n_components=2, random_state=42)
5.     tsne_result = tsne.fit_transform(bow_matrix.toarray())
6.
7.     plt.figure(figsize=(8, 6))
8.     scatter = plt.scatter(tsne_result[:, 0], tsne_result[:, 1], c=clusters, cmap='viridis')
9.     plt.title(title)
10.    plt.legend(*scatter.legend_elements(), title="Clusters")
11.    plt.savefig('training_clusters_plot.png')
12.    plt.show()
```

خروجی:



مرحله ۲-۴-

در این مرحله ابتدا از هر دسته داده های یکتا را جدا کردیم و مینیم تعداد را مشخص کردیم و به تعداد مینیمم از داده ها به صورت رندوم انتخاب کردیم، همچنین اگر داده ای تکراری وجود داشته باشد آن را جایگزین نمی کنیم.

```
1. # Step 2-4: Select a uniform sample from each cluster
2. def select_uniform_sample(bow_matrix, clusters, num_samples):
3.     unique_clusters, cluster_counts = np.unique(clusters, return_counts=True)
4.
5.     selected_samples = []
6.
7.     for cluster in unique_clusters:
8.         cluster_indices = np.where(clusters == cluster)[0]
9.         mi = min(num_samples, cluster_counts[cluster])
10.        selected_indices = np.random.choice(cluster_indices, size=mi, replace=False)
11.        selected_samples.extend(selected_indices)
12.
13.    return selected_samples
```

مرحله ۳ بردار سازی متن

مرحله ۳-۱-

در این مرحله توکنایز را فراخوانی کرده ایم و بر روی داده های تست و آموزشی آن را اعمال کرده ایم:

```
1. # Step 3-1: Convert text sequences to numerical sequences using Tokenizer
2. def tokenize_text_sequences(X_train, X_test):
3.     tokenizer = Tokenizer()
4.     tokenizer.fit_on_texts(X_train)
5.
6.     # Convert text sequences to numerical sequences
7.     train_sequences = tokenizer.texts_to_sequences(X_train)
8.     test_sequences = tokenizer.texts_to_sequences(X_test)
9.
10.    return tokenizer, train_sequences, test_sequences
```

مرحله ۳-۲-

در نهایت ماکسیم گیری انجام شده و به اندازه آن پدینگ اعمال شده:

```
1. # Step 3-2: Pad sequences to have a consistent length
2. def pad_text_sequences(train_sequences, test_sequences):
3.     max_length = max(max(len(seq) for seq in train_sequences), max(len(seq) for seq in
test_sequences))
4.
5.     # Pad sequences
6.     padded_train_sequences = pad_sequences(train_sequences, maxlen=max_length, padding='post')
7.     padded_test_sequences = pad_sequences(test_sequences, maxlen=max_length, padding='post')
8.
9.    return padded_train_sequences, padded_test_sequences, max_length
```

#### مرحله ۴- تعبیه کلمات

در این مرحله طبق کد نمونه فایل مورد نظر دانلود شده و سپس امبدینگ بر روی کلمات اعمال شده تا کلمات به بردار تبدیل شوند.

```
1. # Step 4: Download word embeddings and create an embedding matrix
2. def download_and_create_embedding_matrix(tokenizer):
3.
4.     embedding_file_path = f'embeddings/wiki-news-300d-1M.vec'
5.
6.     if not os.path.isfile(embedding_file_path):
7.         print('Downloading word vectors ... ')
8.         url = f'https://dl.fbaipublicfiles.com/fasttext/vectors-english/wiki-news-300d-1M.vec.zip'
9.         urllib.request.urlretrieve(url, f'wiki-news-300d-1M.vec.zip')
10.
11.        print('Unzipping ... ')
12.        with zipfile.ZipFile(f'wiki-news-300d-1M.vec.zip', 'r') as zip_ref:
13.            zip_ref.extractall('embeddings')
14.            print('done.')
15.
16.            os.remove(f'wiki-news-300d-1M.vec.zip')
17.
18.        embeddings_index = {}
19.        with open(embedding_file_path, encoding='utf-8') as f:
20.            for line in f:
21.                values = line.split()
22.                word = values[0]
23.                coefs = np.asarray(values[1:], dtype='float32')
24.                embeddings_index[word] = coefs
25.
26.        embedding_matrix = np.zeros((len(tokenizer.word_index) + 1, 300))
27.        for word, i in tokenizer.word_index.items():
28.            embedding_vector = embeddings_index.get(word)
29.            if embedding_vector is not None:
30.                embedding_matrix[i] = embedding_vector
31.
32.        return embedding_matrix
```

#### مرحله ۵ ساخت مدل دسته بند

#### مرحله ۵-۱-

در این مرحله مطابق اطلاعات داده شده مدل را آموزش دادیم و هیستوری را جهت نمایش نمودار بازگشت داده ایم:

```
1. # Step 5-1: Build and train a GRU model
2. def build_and_train_gru_model(embedding_matrix, max_length, input_dim):
3.     model = Sequential()
4.     model.add(Embedding(input_dim=input_dim, output_dim=embedding_matrix.shape[1],
```



```

5.         weights=[embedding_matrix], input_length=max_length, trainable=False))
6.     model.add(GRU(units=32, dropout=0.2, recurrent_dropout=0.2))
7.     model.add(Dense(1, activation='sigmoid'))
8.
9.     model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
10.
11.     # Train the model
12.     history = model.fit(padded_train_sequences, y_train.iloc[selected_samples_train],
13.                        epochs=10, batch_size=64, validation_split=0.2)
14.
15.     return model, history

```

خروجی کنسول این مرحله:

```

1. Epoch 1/10
2. 4/4 [=====] - 10s 2s/step - loss: 0.6925 - accuracy: 0.5292 - val_loss:
0.6911 - val_accuracy: 0.5500
3. Epoch 2/10
4. 4/4 [=====] - 6s 2s/step - loss: 0.6901 - accuracy: 0.5542 - val_loss:
0.6897 - val_accuracy: 0.5500
5. Epoch 3/10
6. 4/4 [=====] - 6s 2s/step - loss: 0.6890 - accuracy: 0.5542 - val_loss:
0.6887 - val_accuracy: 0.5500
7. Epoch 4/10
8. 4/4 [=====] - 7s 2s/step - loss: 0.6882 - accuracy: 0.5542 - val_loss:
0.6882 - val_accuracy: 0.5500
9. Epoch 5/10
10. 4/4 [=====] - 10s 3s/step - loss: 0.6874 - accuracy: 0.5542 - val_loss:
0.6882 - val_accuracy: 0.5500
11. Epoch 6/10
12. 4/4 [=====] - 9s 2s/step - loss: 0.6875 - accuracy: 0.5542 - val_loss:
0.6883 - val_accuracy: 0.5500
13. Epoch 7/10
14. 4/4 [=====] - 9s 2s/step - loss: 0.6877 - accuracy: 0.5542 - val_loss:
0.6884 - val_accuracy: 0.5500
15. Epoch 8/10
16. 4/4 [=====] - 9s 2s/step - loss: 0.6872 - accuracy: 0.5542 - val_loss:
0.6883 - val_accuracy: 0.5500
17. Epoch 9/10
18. 4/4 [=====] - 7s 2s/step - loss: 0.6879 - accuracy: 0.5542 - val_loss:
0.6882 - val_accuracy: 0.5500
19. Epoch 10/10
20. 4/4 [=====] - 7s 2s/step - loss: 0.6873 - accuracy: 0.5542 - val_loss:
0.6882 - val_accuracy: 0.5500

```

مرحله ۵-۲-

جهت نمایش نمودار از هیستوری مرحله قبل استفاده شده:

```

1. # Step 5-2: Plot loss and accuracy for the GRU model
2. def plot_metrics(history, title):
3.     plt.figure(figsize=(12, 4))
4.
5.     # Plot training & validation accuracy values
6.     plt.subplot(1, 2, 1)
7.     plt.plot(history.history['accuracy'])
8.     plt.plot(history.history['val_accuracy'])

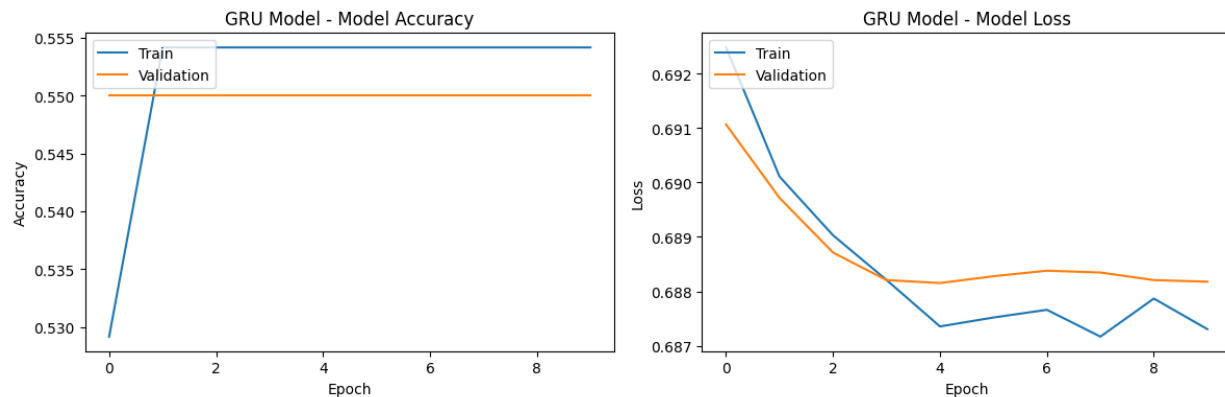
```

```

9.     plt.title(title + ' - Model Accuracy')
10.    plt.xlabel('Epoch')
11.    plt.ylabel('Accuracy')
12.    plt.legend(['Train', 'Validation'], loc='upper left')
13.
14.    # Plot training & validation loss values
15.    plt.subplot(1, 2, 2)
16.    plt.plot(history.history['loss'])
17.    plt.plot(history.history['val_loss'])
18.    plt.title(title + ' - Model Loss')
19.    plt.xlabel('Epoch')
20.    plt.ylabel('Loss')
21.    plt.legend(['Train', 'Validation'], loc='upper left')
22.
23.    plt.tight_layout()
24.    plt.show()
25.

```

خروجی:



مرحله ۵-۳-

در این مرحله نیز مطابق اطلاعات داده شده آموزش مدل انجام شده:

```

1. # Step 5-3: Build and train a BiGRU model
2. def build_and_train_bigru_model(embedding_matrix, max_length, input_dim):
3.     model = Sequential()
4.     model.add(Embedding(input_dim=input_dim, output_dim=embedding_matrix.shape[1],
5.                         weights=[embedding_matrix], input_length=max_length, trainable=False))
6.     model.add(Bidirectional(GRU(units=32, dropout=0.2, recurrent_dropout=0.2)))
7.     model.add(Dense(1, activation='sigmoid'))
8.
9.     model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
10.
11.     # Train the model
12.     history = model.fit(padded_train_sequences, y_train.iloc[selected_samples_train],
13.                       epochs=10, batch_size=64, validation_split=0.2)
14.
15.     return model, history

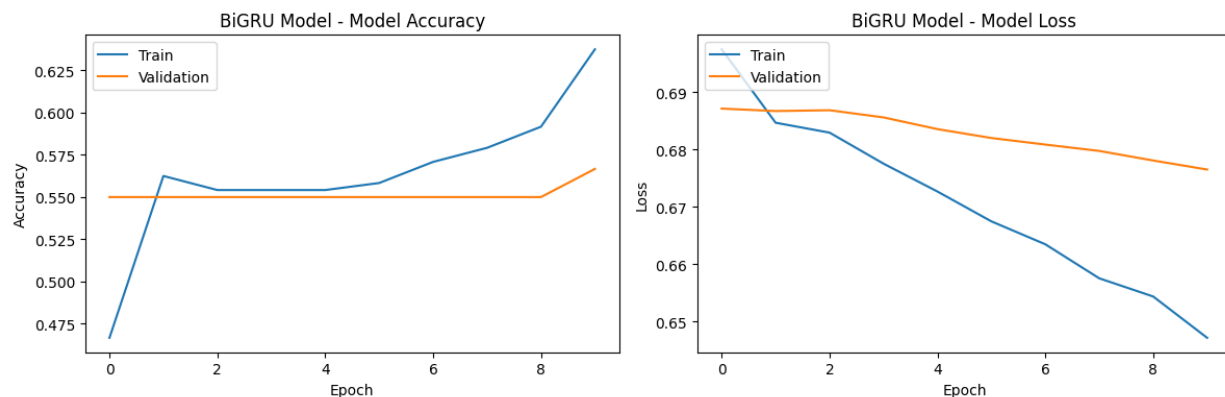
```

خروجی کنسول این مرحله:

```
1. Epoch 1/10
2. 4/4 [=====] - 30s 6s/step - loss: 0.6975 - accuracy: 0.4667 - val_loss:
0.6871 - val_accuracy: 0.5500
3. Epoch 2/10
4. 4/4 [=====] - 22s 5s/step - loss: 0.6846 - accuracy: 0.5625 - val_loss:
0.6867 - val_accuracy: 0.5500
5. Epoch 3/10
6. 4/4 [=====] - 20s 5s/step - loss: 0.6829 - accuracy: 0.5542 - val_loss:
0.6868 - val_accuracy: 0.5500
7. Epoch 4/10
8. 4/4 [=====] - 23s 6s/step - loss: 0.6775 - accuracy: 0.5542 - val_loss:
0.6856 - val_accuracy: 0.5500
9. Epoch 5/10
10. 4/4 [=====] - 23s 6s/step - loss: 0.6726 - accuracy: 0.5542 - val_loss:
0.6835 - val_accuracy: 0.5500
11. Epoch 6/10
12. 4/4 [=====] - 22s 5s/step - loss: 0.6674 - accuracy: 0.5583 - val_loss:
0.6820 - val_accuracy: 0.5500
13. Epoch 7/10
14. 4/4 [=====] - 23s 6s/step - loss: 0.6634 - accuracy: 0.5708 - val_loss:
0.6808 - val_accuracy: 0.5500
15. Epoch 8/10
16. 4/4 [=====] - 22s 5s/step - loss: 0.6575 - accuracy: 0.5792 - val_loss:
0.6797 - val_accuracy: 0.5500
17. Epoch 9/10
18. 4/4 [=====] - 22s 6s/step - loss: 0.6543 - accuracy: 0.5917 - val_loss:
0.6781 - val_accuracy: 0.5500
19. Epoch 10/10
20. 4/4 [=====] - 24s 6s/step - loss: 0.6471 - accuracy: 0.6375 - val_loss:
0.6765 - val_accuracy: 0.5667
```

مرحله ۵-۴-

در این مرحله با استفاده از تابع `plot_metrics` که در مرحله ۵-۳ نوشته شده بود خروجی نموداری گزارش شده:



در نهایت کد های مربوط به اجرای توابع به شرح زیر می باشد:

```
1. file_path = 'IMDB_Dataset.csv'
2. df = load_dataset(file_path)
3.
4. # step1-1
5. convert_labels(df)
6.
7. # step1-2
8. remove_duplicates(df)
9.
10. # step1-3
11. df = random_selection(df)
12.
13. # step1-4
14. df = get_polarity(df)
15.
16. # step1-5
17. plot_sentiment_histograms(df)
18.
19. # step1-6
20. i = 0
21. for index, row in df.iterrows():
22.     print(i)
23.     i += 1
24.     df.at[index, 'review'] = preprocess_text(row['review'])
25.
26. df.to_csv('preprocessed_and_split_IMDB_Dataset.csv', index=False)
27.
28. # step1-7
29.
30. file_path = 'preprocessed_and_split_IMDB_Dataset.csv'
31. df = load_dataset(file_path)
32. X_train, X_test, y_train, y_test = split_dataset(df)
33.
34. # # step2-1
35. bow_matrix_train, bow_matrix_test = create_bow_representation(X_train, X_test)
36.
37. # # Step 2-2
38. train_clusters, test_clusters = kmeans_clustering(bow_matrix_train, bow_matrix_test)
39.
40. # # Step 2-3
41. visualize_clusters(bow_matrix_train, train_clusters, title='Training Set Clusters (t-SNE
Visualization)')
42.
43. # # Step 2-4
44. num_samples_per_cluster = 100 # Adjust the number of samples as needed
45. selected_samples_train = select_uniform_sample(bow_matrix_train, train_clusters,
num_samples_per_cluster)
46. selected_samples_test = select_uniform_sample(bow_matrix_test, test_clusters,
num_samples_per_cluster)
47.
48. # # Step 3-1
49. tokenizer, train_sequences, test_sequences =
tokenize_text_sequences(X_train.iloc[selected_samples_train], X_test.iloc[selected_samples_test])
50.
51. # # Step 3-2
```

```

52. padded_train_sequences, padded_test_sequences, max_length = pad_text_sequences(train_sequences,
test_sequences)
53.
54. # Step 4
55. embedding_matrix = download_and_create_embedding_matrix(tokenizer)
56.
57. # Step 5-1 (GRU Model)
58. gru_model, gru_history = build_and_train_gru_model(embedding_matrix, max_length,
input_dim=len(tokenizer.word_index) + 1)
59.
60. # Step 5-2 (Plot Metrics for GRU Model)
61. plot_metrics(gru_history, 'GRU Model')
62.
63. # Step 5-3 (BiGRU Model)
64. bigru_model, bigru_history = build_and_train_bigru_model(embedding_matrix, max_length,
input_dim=len(tokenizer.word_index) + 1)
65.
66. # Step 5-4 (Plot Metrics for BiGRU Model)
67. plot_metrics(bigru_history, 'BiGRU Model')

```

همچنین دایرکتوری به شرح زیر می باشد:

