

▼ Dependencies

▼ Install

```
1 !pip install contractions
2 !pip install ipython-autotime
3 !pip install fastparquet

Collecting contractions
  Downloading contractions-0.1.73-py2.py3-none-any.whl (8.7 kB)
Collecting textsearch>=0.0.21 (from contractions)
  Downloading textsearch-0.0.24-py2.py3-none-any.whl (7.6 kB)
Collecting anyascii (from textsearch>=0.0.21->contractions)
  Downloading anyascii-0.3.2-py3-none-any.whl (289 kB)
      _____ 289.9/289.9 kB 10.4 MB/s eta 0:00:00
Collecting pyahocorasick (from textsearch>=0.0.21->contractions)
  Downloading pyahocorasick-2.0.0-cp310-cp310-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux2010_x86_64.whl (110 kB)
      _____ 110.8/110.8 kB 16.3 MB/s eta 0:00:00
Installing collected packages: pyahocorasick, anyascii, textsearch, contractions
Successfully installed anyascii-0.3.2 contractions-0.1.73 pyahocorasick-2.0.0 textsearch-0.0.24
Collecting ipython-autotime
  Downloading ipython_autotime-0.3.1-py2.py3-none-any.whl (6.8 kB)
Requirement already satisfied: ipython in /usr/local/lib/python3.10/dist-packages (from ipython-autotime) (7.34.0)
Requirement already satisfied: setuptools>=18.5 in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (67.7.2)
Collecting jedi>=0.16 (from ipython->ipython-autotime)
  Downloading jedi-0.19.1-py2.py3-none-any.whl (1.6 MB)
      _____ 1.6/1.6 MB 21.8 MB/s eta 0:00:00
Requirement already satisfied: decorator in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (4.4.2)
Requirement already satisfied: pickleshare in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (0.7.5)
Requirement already satisfied: traitlets>=4.2 in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (5.7.1)
Requirement already satisfied: prompt-toolkit!=3.0.0,!<3.0.1,<3.1.0,>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (3.0.39)
Requirement already satisfied: pygments in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (2.16.1)
Requirement already satisfied: backcall in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (0.2.0)
Requirement already satisfied: matplotlib-inline in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (0.1.6)
Requirement already satisfied: pexpect>4.3 in /usr/local/lib/python3.10/dist-packages (from ipython->ipython-autotime) (4.8.0)
Requirement already satisfied: parso<0.9.0,>=0.8.3 in /usr/local/lib/python3.10/dist-packages (from jedi>=0.16->ipython->ipython-autotime) (0.8.3)
Requirement already satisfied: ptyprocess>=0.5 in /usr/local/lib/python3.10/dist-packages (from pexpect>4.3->ipython->ipython-autotime) (0.7.0)
Requirement already satisfied: wcwidth in /usr/local/lib/python3.10/dist-packages (from prompt-toolkit!=3.0.0,!<3.0.1,<3.1.0,>=2.0.0->ipython->ipython-autotime) (0.2.8)
Installing collected packages: jedi, ipython-autotime
Successfully installed ipython-autotime-0.3.1 jedi-0.19.1
Collecting fastparquet
  Downloading fastparquet-2023.8.0-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (1.7 MB)
      _____ 1.7/1.7 MB 25.4 MB/s eta 0:00:00
Requirement already satisfied: pandas>=1.5.0 in /usr/local/lib/python3.10/dist-packages (from fastparquet) (1.5.3)
Requirement already satisfied: numpy>=1.20.3 in /usr/local/lib/python3.10/dist-packages (from fastparquet) (1.23.5)
Collecting cramjam>=2.3 (from fastparquet)
  Downloading cramjam-2.7.0-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (1.6 MB)
      _____ 1.6/1.6 MB 103.8 MB/s eta 0:00:00
Requirement already satisfied: fsspec in /usr/local/lib/python3.10/dist-packages (from fastparquet) (2023.6.0)
Requirement already satisfied: packaging in /usr/local/lib/python3.10/dist-packages (from fastparquet) (23.2)
Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.5.0->fastparquet) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=1.5.0->fastparquet) (2023.3.post1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.1->pandas>=1.5.0->fastparquet) (1.16.0)
Installing collected packages: cramjam, fastparquet
Successfully installed cramjam-2.7.0 fastparquet-2023.8.0
```

▼ Imports

```
1 import os
2 import re
3 import shutil
4 import unicodedata
5 import multiprocessing
6
7 import warnings
8 warnings.filterwarnings("ignore")
9
10 import numpy as np
11 import pandas as pd
12 import requests
13
14 import nltk
15 from nltk.corpus import stopwords, wordnet
16 from nltk.stem import WordNetLemmatizer
17 from nltk.tokenize import word_tokenize
18
19 nltk.download('punkt', quiet=True)
20 nltk.download('wordnet', quiet=True)
21 nltk.download('stopwords', quiet=True)
22 nltk.download('averaged_perceptron_tagger', quiet=True)
23
24 import contractions
25
26 import gensim
27 import gensim.downloader as api
28 from gensim.models import Word2Vec
29
30 from sklearn.model_selection import train_test_split
31 from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score
32
33 from sklearn.linear_model import Perceptron
34 from sklearn.svm import LinearSVC
35
36 import torch
37 import torch.nn as nn
38 import torch.optim as optim
39 from torch.utils.data.sampler import RandomSampler, BatchSampler
40 from torch.utils.data import Dataset, DataLoader
41
42 from tqdm.notebook import tqdm
43
44 %load_ext autotime


time: 309 μs (started: 2023-10-20 19:58:20 +00:00)
```

▼ Config

Set up important configuration parameters and file paths for the project, making it easy to manage various settings and paths from one centralized location

Place the `amazon_reviews_us_Office_Products_v1_00.tsv.gz` at the same level as `noetbook`

```
1 # os.chdir("/content/drive/MyDrive/Colab Notebooks/CSCI544/HW3")
2 os.environ['CUDA_LAUNCH_BLOCKING'] = "1"
3
4 CURRENT_DIR = os.getcwd()
5
6
7 class DatasetConfig:
8     RANDOM_STATE = 34
9     TEST_SPLIT = 0.2
10    N_SAMPLES_EACH_CLASS = 50000
11    DATA_PATH = os.path.join(
12        CURRENT_DIR, "amazon_reviews_us_Office_Products_v1_00.tsv.gz"
13    )
14    PROCESSED_DATA_PATH = os.path.join(
15        CURRENT_DIR, "amazon_review_processed_sentiment_analysis.parquet"
16    )
17    PREPROCESSED_DATA_PATH = os.path.join(
18        CURRENT_DIR, "amazon_review_preprocessed_sentiment_analysis.parquet"
19    )
20    BUILD_NEW = True
21    if os.path.exists(PROCESSED_DATA_PATH) and os.path.exists(PREPROCESSED_DATA_PATH):
22        BUILD_NEW = False
23
24
25 class Word2VecConfig:
26     PRETRAINED_MODEL = "word2vec-google-news-300"
27     PRETRAINED_DEFAULT_SAVE_PATH = os.path.join(
28         gensim.downloader.BASE_DIR, PRETRAINED_MODEL, f"{PRETRAINED_MODEL}.gz"
29     )
30     PRETRAINED_MODEL_SAVE_PATH = os.path.join(
31         CURRENT_DIR, PRETRAINED_MODEL, f"{PRETRAINED_MODEL}.gz"
32     )
33     WINDOW_SIZE = 13
34     MAX_LENGTH = 300
35     EMBEDDING_SIZE = 300
36     MIN_WORD_COUNT = 9
37     CUSTOM_MODEL_PATH = os.path.join(CURRENT_DIR, "word2vec-custom.model")
```

 time: 2.43 ms (started: 2023-10-20 22:01:48 +00:00)

▼ Helper Functions

▼ Download & Save Pretrained model

- Run the `api.load()` once and copied the model from temporary path to local drive for fast loading of model in memory.

References:

1. [Faster way to load word2vec model](#)
2. [Tutorial](#)

```
1 def load_pretrained_model():
2     if not os.path.exists(Word2VecConfig.PRETRAINED_MODEL_SAVE_PATH):
3         # Create a directory if it doesn't exist
4         os.makedirs(Word2VecConfig.PRETRAINED_MODEL, exist_ok=True)
5         # Download the model embeddings
6         pretrained_model = api.load(Word2VecConfig.PRETRAINED_MODEL, return_path=True)
7         # Copy & save the embeddings file
8         shutil.copyfile(
9             Word2VecConfig.PRETRAINED_DEFAULT_SAVE_PATH, Word2VecConfig.PRETRAINED_MODEL
10        )
11    else:
12        pretrained_model = gensim.models.keyedvectors.KeyedVectors.load_word2vec_format(
13            Word2VecConfig.PRETRAINED_MODEL_SAVE_PATH, binary=True
14        )
15    return pretrained_model
16
17
18 # Load the pretrained model
19 pretrained_model = load_pretrained_model()
```

time: 1min 31s (started: 2023-10-20 06:45:23 +00:00)

▼ Accelerator Configuration

```
1 def get_device():
2     if torch.cuda.is_available():
3         # Check if GPU is available
4         return torch.device("cuda")
5     else:
6         # Use CPU if no GPU or TPU is available
7         return torch.device("cpu")
8
9 device = get_device()
10 device
```

device(type='cpu')time: 5.32 ms (started: 2023-10-20 19:58:39 +00:00)

▼ Download Data

Checks if a file specified by `DatasetConfig.DATA_PATH` exists. If not, it downloads the file from a given URL and saves it with the same name. If the file already exists, it prints a message indicating so

```
1 if not os.path.exists(DatasetConfig.DATA_PATH):
2     url = (
3         "https://web.archive.org/web/20201127142707if_/https://s3.amazonaws.com/amazon-reviews-pds"
4         "/tsv/amazon_reviews_us_Office_Products_v1_00.tsv.gz"
5     )
6     file_name = DatasetConfig.DATA_PATH
7
8     # Stream and download heavy file in chunks
9     with requests.get(url, stream=True) as response:
10         with open(file_name, "wb") as file:
11             for chunk in response.iter_content(chunk_size=8192):
12                 file.write(chunk)
13
14     print(f"Downloaded file '{os.path.relpath(file_name)}' successfully.")
```

```
15 else:
    File '/content/drive/MyDrive/Colab Notebooks/CSCI544/HW3/amazon_reviews_us_Office_Products_v1_00.tsv.gz' already exists.
    time: 1.93 ms (started: 2023-10-20 19:58:42 +00:00)
```

▼ Dataset Preparation

This code provides a pipeline for processing and preparing a dataset for sentiment analysis:

1. `LoadData` class loads a dataset from a specified path, keeping only relevant columns.
2. `ProcessData` class performs the following tasks:
 - Converts star ratings to numeric values.
 - Classifies sentiments based on star ratings (1 for negative, 2 for positive).
 - Balances the dataset by sampling an equal number of samples for both sentiments.
3. `CleanText` class defines various text cleaning operations:
 - Removing non-ASCII characters.
 - Expanding contractions.
 - Removing email addresses, URLs, and HTML tags.
 - Lowercasing and stripping spaces.
4. `clean_and_process_data` function executes the entire data processing pipeline:
 - Loads the data.
 - Applies basic processing.
 - Balances the dataset.
 - Cleans the text.
 - Tokenizes the reviews.
5. `preprocess_review_body` function generates word embeddings for each word in a review using a pre-trained Word2Vec model.
6. `get_reviews_dataset` function handles the entire data preprocessing and embedding generation process. It checks if the preprocessed data already exists, and if not, it performs the data preprocessing and saves the preprocessed data in Parquet format.

Overall, this pipeline ensures that the dataset is properly loaded, cleaned, processed, balanced, and transformed into embeddings suitable for sentiment analysis.

- Note:
- Parquet format is efficient for storage.
 - Storing data to avoid running the pipeline and embedding generation process all over again.
 - Provides a ready-to-use dataset for sentiment analysis tasks, allowing for quicker experimentation and model training

▼ Read and Process

```
1 class LoadData:
2     @staticmethod
3     def load_data(path):
4         df = pd.read_csv(
5             path,
6             sep="\t",
7             usecols=["review_headline", "review_body", "star_rating"],
8             on_bad_lines="skip",
9             memory_map=True,
10        )
11        return df
12
13
14 class ProcessData:
15     @staticmethod
16     def filter_columns(df):
17         return df.loc[:, ["review_body", "star_rating"]]
18
19     @staticmethod
20     def convert_star_rating(df):
21         df["star_rating"] = pd.to_numeric(df["star_rating"], errors="coerce")
22         df.dropna(subset=["star_rating"], inplace=True)
23         return df
24
25     @staticmethod
26     def classify_sentiment(df):
27         df["sentiment"] = df["star_rating"].apply(lambda x: 1 if x <= 3 else 2)
28         return df
29
30     @staticmethod
31     def sample_data(df, n_samples, random_state):
32         sampled_df = pd.concat(
33             [
34                 df.query("sentiment==1").sample(n=n_samples, random_state=random_state),
35                 df.query("sentiment==2").sample(n=n_samples, random_state=random_state),
36             ],
37             ignore_index=True,
38         ).sample(frac=1, random_state=random_state, ignore_index=True)
39
40         sampled_df.drop(columns=["star_rating"], inplace=True)
41         return sampled_df
42
43
44 class CleanText:
45     @staticmethod
46     def unicode_to_ascii(s):
47         return "".join(
48             c for c in unicodedata.normalize("NFD", s) if unicodedata.category(c) != "Mn"
49         )
50
51     @staticmethod
52     def expand_contractions(text):
53         """Expand contraction for eg., wouldn't => would not"""
54         return contractions.fix(text)
55
56     @staticmethod
57     def remove_email_addresses(text):
58         return re.sub(r"[a-zA-Z0-9_\-\.]+@[a-zA-Z0-9_\-\.]+\.[a-zA-Z]{2,5}", "", text)
59
60     @staticmethod
61     def remove_urls(text):
62         return re.sub(r"\bhttps?:\/\/\S+|www\.\S+", "", text)
63
64     @staticmethod
```

```
65 def remove_html_tags(text):
66     return re.sub(r"<.*?>", "", text)
67
68 @staticmethod
69 def clean_text(text):
70     text = text.lower().strip()
71     text = CleanText.unicode_to_ascii(text)
72     # text = CleanText.remove_email_addresses(text)
73     # text = CleanText.remove_urls(text)
74     text = CleanText.remove_html_tags(text)
75     text = CleanText.expand_contractions(text)
76
77     # creating a space between a word and the punctuation following it
78     # text = re.sub(r"([?.!,;])", r" \1 ", text)
79     # text = re.sub(r'[" "]+', " ", text)
80
81     # removes all non-alphabetical characters
82     # text = re.sub(r"[^a-zA-Z\s]+", "", text)
83
84     # remove extra spaces
85     # text = re.sub(" +", " ", text)
86     return text
87
88
89 def clean_and_process_data(path):
90     df = LoadData.load_data(path)
91
92     # Basic processing
93     df_filtered = ProcessData.filter_columns(df)
94     df_filtered = ProcessData.convert_star_rating(df_filtered)
95     df_filtered = ProcessData.classify_sentiment(df_filtered)
96
97     balanced_df = ProcessData.sample_data(
98         df_filtered, DatasetConfig.N_SAMPLES_EACH_CLASS, DatasetConfig.RANDOM_STATE
99     )
100
101     # Clean data
102     balanced_df.dropna(inplace=True)
103     balanced_df["review_body"] = balanced_df["review_body"].astype(str)
104     balanced_df["review_body"] = balanced_df["review_body"].apply(CleanText.clean_text)
105     # Drop reviews that are empty
106     balanced_df = balanced_df.loc[balanced_df["review_body"].str.strip() != ""]
107
108     # Tokenize Reviews
109     balanced_df["review_body"] = balanced_df["review_body"].apply(word_tokenize)
110     return balanced_df
111
112
113 def preprocess_review_body(text, word2vec_model, topn=None):
114     embeddings = [word2vec_model[word] for word in text if word in word2vec_model]
115
116     if topn is not None:
117         embeddings = np.concatenate(embeddings[:topn], axis=0)
118     else:
119         embeddings = np.mean(embeddings, axis=0)
120     return embeddings
121
122
123 def get_reviews_dataset(new=False):
124     if new or not os.path.exists(DatasetConfig.DATA_PATH):
125         balanced_df = clean_and_process_data(DatasetConfig.DATA_PATH)
126         balanced_df.to_parquet(DatasetConfig.PROCESSED_DATA_PATH, index=False)
127
128         # Preprocess data and generate word2vec embeddings Avg and top 10
129         balanced_df["embeddings"] = balanced_df["review_body"].apply(
130             lambda text: preprocess_review_body(text, pretrained_model, topn=None)
131         )
132         # Drop rows with NaN embeddings
133         balanced_df.dropna(subset=["embeddings"], inplace=True)
134
135         balanced_df["embeddings_top_10"] = balanced_df["review_body"].apply(
136             lambda text: preprocess_review_body(text, pretrained_model, topn=10)
137         )
138
139         balanced_df.to_parquet(DatasetConfig.PREPROCESSED_DATA_PATH, index=False)
140     else:
141         balanced_df = pd.read_parquet(
142             DatasetConfig.PREPROCESSED_DATA_PATH,
143             # engine="fastparquet"
144         )
145     return balanced_df
146
147
148 time: 2.38 ms (started: 2023-10-20 19:59:52 +00:00)
```

```
1 balanced_df = get_reviews_dataset(
2     new=DatasetConfig.BUILD_NEW
3 )
4 print("Total Records:", balanced_df.shape)
5 balanced_df.head(10)
```

Total Records: (99862, 4)				
	review_body	sentiment	embeddings	embeddings_top_10
0	[i, set, up, a, photo, booth, at, my, sister, ...	2	[0.016994974, 0.024544675, -0.010975713, 0.093...	[-0.22558594, -0.01953125, 0.09082031, 0.23730...
1	[like, everyone, else, ,, i, like, saving, mon...	1	[0.044110615, 0.036876563, 0.0371785, 0.113560...	[0.103515625, 0.13769531, -0.0029754639, 0.181...
2	[the, pen, is, perfect, what, i, want, !, howe...	2	[0.026102701, 0.029064532, 0.010800962, 0.0622...	[0.080078125, 0.10498047, 0.049804688, 0.05346...
3	[i, think, they, are, too, expensive, for, the...	1	[-0.0039075767, 0.032967318, 0.02339106, 0.113...	[-0.22558594, -0.01953125, 0.09082031, 0.23730...
4	[black, is, working, wonderfully, ,, and, both...	1	[0.034285888, 0.013478661, 0.041618653, 0.1132...	[0.10498047, 0.018432617, 0.008972168, -0.0128...
5	[i, have, problems, with, the, moveable, tab, ...	1	[0.010405041, 0.026173819, 0.03433373, 0.09698...	[-0.22558594, -0.01953125, 0.09082031, 0.23730...
6	[this, printer, sucks, !, it, started, out, wo...	1	[0.05581854, 0.035414256, 0.047512088, 0.09278...	[0.109375, 0.140625, -0.03173828, 0.16601562, ...

▼ Review Body stats

- Mean number of words = 66
- Median number of words = 37
- Limiting sequence length for RNN based embeddings = 45

```
1 balanced_df["review_body"].apply(len).describe().round(2)

count      99862.00
mean         65.94
std         100.17
min           1.00
25%          20.00
50%          37.00
75%          76.00
max         4847.00
Name: review_body, dtype: float64time: 272 ms (started: 2023-10-20 20:00:20 +00:00)
```

▼ Train, Valid and Test Spilts

```
1 # Create train and temp sets (80% train, 20% valid + test)
2 train_df, valid_df = train_test_split(
3     balanced_df,
4     test_size=0.20,
5     random_state=DatasetConfig.RANDOM_STATE,
6     stratify=balanced_df["sentiment"]
7 )
8
9 # Create valid and test sets (15% valid, 5% test)
10 valid_df, test_df = train_test_split(
11     valid_df,
12     test_size=0.25, # 25% of 20% is 5%
13     random_state=DatasetConfig.RANDOM_STATE,
14     stratify=valid_df["sentiment"]
15 )
```

time: 83.1 ms (started: 2023-10-20 20:04:14 +00:00)

▼ Word Embedding

▼ Semantic similarity examples with pretrained embeddings

```
1 # Example 1: King - Man + Woman = Queen
2 result = pretrained_model.most_similar(positive=['woman', 'king'], negative=['man'])
3 print(f"Semantic Similarity: {result[0][0]}")
4
5 # Example 2: excellent ~ outstanding
6 result = pretrained_model.similarity('excellent', 'outstanding')
7 print(f"Semantic Similarity: {result}")
8
9 # Example 3: Paris - France + Italy = Milan
10 result = pretrained_model.most_similar(positive=['Italy', 'Paris'], negative=['France'])
11 print(f"Semantic Similarity: {result[0][0]}")
12
13 # Example 4: Car - Wheel + Boat = Yacht
14 result = pretrained_model.most_similar(positive=['Boat', 'Car'], negative=['Wheel'])
15 print(f"Semantic Similarity: {result[0][0]}")
16
17 # Example 5: Delicious ~ Tasty
18 result = pretrained_model.similarity('Delicious', 'Tasty')
19 print(f"Semantic Similarity: {result}")
20
21 # Example 6: Computer ~ Plant
22 result = pretrained_model.similarity('Computer', 'Plant')
23 print(f"Semantic Similarity: {result}")
24
25 # Example 7: Cat ~ Dog
26 result = pretrained_model.similarity('Cat', 'Dog')
27 print(f"Semantic Similarity: {result}")
```

Semantic Similarity: queen
Semantic Similarity: 0.5567485690116882
Semantic Similarity: Milan
Semantic Similarity: Yacht
Semantic Similarity: 0.5718502402305603
Semantic Similarity: 0.04445184767246246
Semantic Similarity: 0.6061107516288757
time: 9.78 s (started: 2023-10-18 21:26:10 +00:00)

```
1 del pretrained_model
```

time: 479 µs (started: 2023-10-20 06:47:24 +00:00)

▼ Custom Word2Vec Embeddings Generation

```
1 sentences=train_df["review_body"].apply(lambda x: x.tolist()).tolist()
2
3 # Train Word2Vec model
4 w2v_model_custom = Word2Vec(
5     sentences=sentences,
6     vector_size=Word2VecConfig.MAX_LENGTH,
7     window=Word2VecConfig.WINDOW_SIZE,
8     min_count=Word2VecConfig.MIN_WORD_COUNT,
9     workers=multiprocessing.cpu_count()
10 )
11
12 # Save the model
13 w2v_model_custom.save(Word2VecConfig.CUSTOM_MODEL_PATH)
```

time: 1min 30s (started: 2023-10-18 21:36:09 +00:00)

▼ Test Custom Embeddings

```
1 # Load the custom model
2 w2v_model_custom = Word2Vec.load(Word2VecConfig.CUSTOM_MODEL_PATH)
3
4 # Example 1: King - Man + Woman = Queen
5 res = w2v_model_custom.wv.most_similar(positive=['woman', 'king'], negative=['man'])
6 print(f"Semantic Similarity (Custom Model): {res[0]}")
7
8 # Example 2: excellent ~ outstanding
9 res = w2v_model_custom.wv.similarity('excellent', 'outstanding')
10 print(f"Semantic Similarity (Custom Model): {res}")
```

Semantic Similarity (Custom Model): ('queen', 0.5723455548286438)
Semantic Similarity (Custom Model): 0.7957370281219482
time: 241 ms (started: 2023-10-18 21:37:47 +00:00)

▼ Conclusion

What do you conclude from comparing vectors generated by yourself and the pretrained model? Which of the Word2Vec models seems to encode semantic similarities between words better?

1. Custom-trained Word2Vec Model:
 - Strengths:
 - Captures domain-specific relationships and nuances as it trained on very specific dataset.
 - Weaknesses:
 - It may not perform as well on tasks outside of its training domain.
 - The quality of embeddings heavily depends on the dataset used for training.
 - For example, if the dataset is small or not representative of the overall language, the embeddings may be less reliable.
2. Pretrained "word2vec-google-news-300" Model:
 - Strengths:
 - This model has been pretrained on a massive corpus of text from various domains, making it highly versatile and capable of capturing a wide range of semantic relationships.
 - It can generalize well to different tasks and domains.
 - Weaknesses:
 - While it provides strong generalization, it may not capture domain-specific relationships as effectively as a model trained on domain-specific data.
- The semantic similarity score is higher for the pretrained model compared to the custom model. This indicates that the pretrained model is better at encoding semantic similarities between words.
 - The custom Word2Vec model, which was trained on the provided dataset, may not have had access to as diverse and extensive a corpus as the pretrained model. This can lead to limitations in its ability to generalize and capture nuanced semantic relationships.

```
1 del w2v_model_custom, res, sentences

time: 319 µs (started: 2023-10-19 01:33:10 +00:00)
```

▼ Simple Models

```
1 def evaluate_model(model, X_test, y_test):
2     # Predict on the test set
3     y_pred = model.predict(X_test)
4
5     # Calculate evaluation metrics
6     precision = precision_score(y_test, y_pred, average="binary")
7     recall = recall_score(y_test, y_pred, average="binary")
8     f1 = f1_score(y_test, y_pred, average="binary")
9     accuracy = accuracy_score(y_test, y_pred)
10
11     return precision, recall, f1, accuracy
12
13
14 def train_and_evaluate_model(model_class, X_train, y_train, X_test, y_test, **model_params):
15     # Initialize model
16     model = model_class(**model_params)
17
18     # Train the model
19     model.fit(X_train, y_train)
20
21     # Evaluate model
22     precision, recall, f1, accuracy = evaluate_model(model, X_test, y_test)
23     return model, precision, recall, f1, accuracy

time: 1.21 ms (started: 2023-10-20 06:47:33 +00:00)
```

```
1 X_train = np.vstack(train_df["embeddings"])
2 y_train = train_df["sentiment"]
3 X_test = np.vstack(test_df["embeddings"])
4 y_test = test_df["sentiment"]

time: 338 ms (started: 2023-10-20 20:04:40 +00:00)
```

▼ SVM

Params	Precision	Recall	F1	Accuracy	Features Used
LinearSVC(C=0.1, max_iter=10000)	0.7997	0.8671	0.8320	0.8321	Word2Vec
LinearSVC(max_iter=10000)	0.8045	0.8623	0.8324	0.8262	Word2Vec
LinearSVC(C=0.01, max_iter=15000)	0.7836	0.8835	0.8305	0.8281	Word2Vec

```
1 # Train and evaluate LinearSVC model
2 (
3     _,
4     precision_svc,
5     recall_svc,
6     f1_svc,
7     acc_svc
8 ) = train_and_evaluate_model(
9     LinearSVC,
10    X_train, y_train, X_test, y_test,
11    max_iter=10000,
12    # C=0.1
13 )
14
15 print(f'Precision Recall F1 Accuracy (LinearSVC): {precision_svc:.4f} {recall_svc:.4f} {f1_svc:.4f} {acc_svc:.4f}')
```

Precision Recall F1 Accuracy (LinearSVC): 0.8045 0.8623 0.8324 0.8262
time: 28.2 s (started: 2023-10-19 01:31:55 +00:00)

▼ Perceptron

Params	Precision	Recall	F1	Accuracy	Features Used
Perceptron(eta0=0.01, max_iter=5000, penalty='elasticnet', warm_start=True)	0.7693	0.8778	0.8200	0.8071	Word2Vec

Params	Precision	Recall	F1	Accuracy	Features Used
Perceptron(max_iter=5000)	0.7786	0.8613	0.8179	0.8110	Word2Vec
Perceptron()	0.7786	0.8613	0.8179	0.8110	Word2Vec
Perceptron(eta0=0.1, max_iter=5000, penalty='elasticnet', warm_start=True)	0.5977	0.9844	0.7438	0.6655	Word2Vec
Perceptron(eta0=0.001, max_iter=10000, penalty='l2')	0.7367	0.9114	0.8148	0.7849	Word2Vec
Perceptron(eta0=0.01, max_iter=10000, penalty='l2', warm_start=True)	0.7653	0.8789	0.8181	0.8002	Word2Vec
Perceptron(eta0=0.01, penalty='l1', warm_start=True)	0.6133	0.9813	0.7548	0.6832	Word2Vec

```
1 # Train and evaluate Perceptron model using Bow features
2 (
3     _,
4     precision_perceptron,
5     recall_perceptron,
6     f1_perceptron,
7     acc_perceptron
8 ) = train_and_evaluate_model(
9     Perceptron,
10    X_train, y_train, X_test, y_test,
11    max_iter=5000,
12    eta0=0.01,
13    warm_start=True,
14    penalty="elasticnet"
15 )
16
17 print(f'Precision Recall F1 (Perceptron): {precision_perceptron:.4f} {recall_perceptron:.4f} {f1_perceptron:.4f} {acc_perceptron:.4f}')
```

Precision Recall F1 (Perceptron): 0.7693 0.8778 0.8200 0.8071
time: 1.1 s (started: 2023-10-19 01:32:24 +00:00)

▼ With TFIDF Features

▼ Homework 1 Script Edited

```
1  # @title Homework 1 Script Edited
2
3  %%writefile HW1-CSCI544-wo-neg-sw.py
4  # Python Version: 3.10.12
5
6  import re
7  import unicodedata
8
9  import warnings
10
11  warnings.filterwarnings("ignore")
12
13  import numpy as np
14  import pandas as pd
15
16  import nltk
17  from nltk.corpus import stopwords, wordnet
18  from nltk.stem import WordNetLemmatizer
19  from nltk.tokenize import word_tokenize
20
21  nltk.download("punkt", quiet=True)
22  nltk.download("wordnet", quiet=True)
23  nltk.download("stopwords", quiet=True)
24  nltk.download("averaged_perceptron_tagger", quiet=True)
25
26  import contractions
27
28  from sklearn.model_selection import train_test_split
29  from sklearn.feature_extraction.text import TfidfVectorizer
30  from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score
31
32  from sklearn.linear_model import Perceptron
33  from sklearn.svm import LinearSVC
34
35
36  class Config:
37      RANDOM_STATE = 56
38      DATA_PATH = "amazon_reviews_us_Office_Products_v1_00.tsv.gz"
39      TEST_SPLIT = 0.2
40      N_SAMPLES_EACH_CLASS = 50000
41      NUM_TFIDF_FEATURES = 5000
42      NUM_BOW_FEATURES = 5000
43
44
45  class DataLoader:
46      @staticmethod
47      def load_data(path):
48          df = pd.read_csv(
49              path,
50              sep="\t",
51              usecols=["review_headline", "review_body", "star_rating"],
52              on_bad_lines="skip",
53              memory_map=True,
54          )
55          return df
56
57
58  class DataProcessor:
59      @staticmethod
60      def filter_columns(df):
61          return df.loc[:, ["review_body", "star_rating"]]
62
63      @staticmethod
64      def convert_star_rating(df):
65          df["star_rating"] = pd.to_numeric(df["star_rating"], errors="coerce")
66          df.dropna(subset=["star_rating"], inplace=True)
67          return df
68
69      @staticmethod
70      def classify_sentiment(df):
71          df["sentiment"] = df["star_rating"].apply(lambda x: 1 if x <= 3 else 2)
72          return df
73
74      @staticmethod
75      def sample_data(df, n_samples, random_state):
76          sampled_df = pd.concat(
77              [
78                  df.query("sentiment==1").sample(n=n_samples, random_state=random_state),
79                  df.query("sentiment==2").sample(n=n_samples, random_state=random_state),
80              ],
81              ignore_index=True,
```

```

82         ).sample(frac=1, random_state=random_state)
83
84         sampled_df.drop(columns=["star_rating"], inplace=True)
85         return sampled_df
86
87
88     class TextCleaner:
89         @staticmethod
90         def unicode_to_ascii(s):
91             return "".join(
92                 c for c in unicodedata.normalize("NFD", s) if unicodedata.category(c) != "Mn"
93             )
94
95         @staticmethod
96         def expand_contractions(text):
97             return contractions.fix(text)
98
99         @staticmethod
100        def remove_email_addresses(text):
101            return re.sub(r"[a-zA-Z0-9_\-\.]+@[a-zA-Z0-9_\-\.]+\.[a-zA-Z]{2,5}", " ", text)
102
103        @staticmethod
104        def remove_urls(text):
105            return re.sub(r"\bhttps?:\/\/\S+|www\.\S+", " ", text)
106
107        @staticmethod
108        def remove_html_tags(text):
109            return re.sub(r"<.*>", "", text)
110
111        @staticmethod
112        def clean_text(text):
113            text = TextCleaner.unicode_to_ascii(text.lower().strip())
114            # replacing email addresses with empty string
115            text = TextCleaner.remove_email_addresses(text)
116            # replacing urls with empty string
117            text = TextCleaner.remove_urls(text)
118            # Remove HTML tags
119            text = TextCleaner.remove_html_tags(text)
120            # Expand contraction for eg., wouldn't => would not
121            text = TextCleaner.expand_contractions(text)
122            # creating a space between a word and the punctuation following it
123            text = re.sub(r"([?.!,;])", r" \1 ", text)
124            text = re.sub(r'[" "]', " ", text)
125            # removes all non-alphabetical characters
126            text = re.sub(r"[^a-zA-Z\s]", "", text)
127            # remove extra spaces
128            text = re.sub(" +", " ", text)
129            text = text.strip()
130            return text
131
132
133        class TextPreprocessor:
134            lemmatizer = WordNetLemmatizer()
135
136            @staticmethod
137            def get_stopwords_pattern():
138                # Stopword list
139                og_stopwords = set(stopwords.words("english"))
140
141                # Define a list of negative words to remove
142                neg_words = ["no", "not", "nor", "neither", "none", "never", "nobody", "nowhere"]
143                custom_stopwords = [word for word in og_stopwords if word not in neg_words]
144                pattern = re.compile(r"\b(" + r"|".join(custom_stopwords) + r")\b\s*")
145                return pattern
146
147            @staticmethod
148            def pos_tagger(tag):
149                if tag.startswith("J"):
150                    return wordnet.ADJ
151                elif tag.startswith("V"):
152                    return wordnet.VERB
153                elif tag.startswith("N"):
154                    return wordnet.NOUN
155                elif tag.startswith("R"):
156                    return wordnet.ADV
157                else:
158                    return None
159
160            @staticmethod
161            def lemmatize_text_using_pos_tags(text):
162                words = nltk.pos_tag(word_tokenize(text))
163                words = map(lambda x: (x[0], TextPreprocessor.pos_tagger(x[1])), words)
164                lemmatized_words = [
165                    TextPreprocessor.lemmatizer.lemmatize(word, tag) if tag else word for word, tag in words
166                ]
167                return " ".join(lemmatized_words)
168
169            @staticmethod
170            def lemmatize_text(text):
171                words = word_tokenize(text)
172                lemmatized_words = [TextPreprocessor.lemmatizer.lemmatize(word) for word in words]
173                return " ".join(lemmatized_words)
174
175            pattern = get_stopwords_pattern()
176
177            @staticmethod
178            def preprocess_text(text):
179                # replacing all the stopwords
180                text = TextPreprocessor.pattern.sub("", text)
181                text = TextPreprocessor.lemmatize_text(text)
182                return text
183
184
185        clean_text_vect = np.vectorize(TextCleaner.clean_text)
186        preprocess_text_vect = np.vectorize(TextPreprocessor.preprocess_text)
187
188
189        def clean_and_process_data(path):
190            df = DataLoader.load_data(path)
191            df_filtered = DataProcessor.filter_columns(df)
192            df_filtered = DataProcessor.convert_star_rating(df_filtered)
193            df_filtered = DataProcessor.classify_sentiment(df_filtered)
194
195            balanced_df = DataProcessor.sample_data(
196                df_filtered, Config.N_SAMPLES_EACH_CLASS, Config.RANDOM_STATE
197            )
198
199            balanced_df["review_body"] = balanced_df["review_body"].astype(str)
200
201            # Clean data

```



```
201     # Clean data
202     # avg_len_before_clean = balanced_df["review_body"].apply(len).mean()
203     balanced_df["review_body"] = balanced_df["review_body"].apply(clean_text_vect)
204     # Drop reviews that are empty
205     balanced_df = balanced_df.loc[balanced_df["review_body"].str.strip() != ""]
206     # avg_len_after_clean = balanced_df["review_body"].apply(len).mean()
207
208     # Preprocess data
209     # avg_len_before_preprocess = avg_len_after_clean
210     balanced_df["review_body"] = balanced_df["review_body"].apply(preprocess_text_vect)
211     # avg_len_after_preprocess = balanced_df["review_body"].apply(len).mean()
212
213     # Print Results
214     # print(f"{avg_len_before_clean:.2f}, {avg_len_after_clean:.2f}")
215     # print(f"{avg_len_before_preprocess:.2f}, {avg_len_after_preprocess:.2f}")
216
217     return balanced_df
218
219
220 def evaluate_model(model, X_test, y_test):
221     # Predict on the test set
222     y_pred = model.predict(X_test)
223
224     # Calculate evaluation metrics
225     precision = precision_score(y_test, y_pred, average="binary")
226     recall = recall_score(y_test, y_pred, average="binary")
227     f1 = f1_score(y_test, y_pred, average="binary")
228     accuracy = accuracy_score(y_test, y_pred)
229
230     return precision, recall, f1, accuracy
231
232
233 def train_and_evaluate_model(model_class, X_train, y_train, X_test, y_test, **model_params):
234     # Initialize model
235     model = model_class(**model_params)
236
237     # Train the model
238     model.fit(X_train, y_train)
239
240     # Evaluate model
241     precision, recall, f1, accuracy = evaluate_model(model, X_test, y_test)
242     return model, precision, recall, f1, accuracy
243
244
245 def main():
246     balanced_df = clean_and_process_data(Config.DATA_PATH)
247
248     # Splitting the reviews dataset
249     X_train, X_test, y_train, y_test = train_test_split(
250         balanced_df["review_body"],
251         balanced_df["sentiment"],
252         test_size=Config.TEST_SPLIT,
253         random_state=Config.RANDOM_STATE,
254     )
255
256     # Feature Extraction
257     tfidf_vectorizer = TfidfVectorizer(max_features=Config.NUM_TFIDF_FEATURES)
258     X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)
259     X_test_tfidf = tfidf_vectorizer.transform(X_test)
260
261     # Train and evaluate Perceptron model using TF-IDF features
262     (
263         _,
264         precision_perceptron_tfidf,
265         recall_perceptron_tfidf,
266         f1_perceptron_tfidf,
267         acc_perceptron_tfidf
268     ) = train_and_evaluate_model(
269         Perceptron, X_train_tfidf, y_train, X_test_tfidf, y_test, max_iter=4000
270     )
271
272     # Train and evaluate SVM model using TF-IDF features
273     (
274         _,
275         precision_svm_tfidf,
276         recall_svm_tfidf,
277         f1_svm_tfidf,
278         acc_svm_tfidf
279     ) = train_and_evaluate_model(
280         LinearSVC, X_train_tfidf, y_train, X_test_tfidf, y_test, max_iter=2500
281     )
282
283     # Print the results
284     print("Precision Recall F1-Score Accuracy")
285     print("Perceptron")
286     print(
287         f"{precision_perceptron_tfidf:.4f} {recall_perceptron_tfidf:.4f} {f1_perceptron_tfidf:.4f} {acc_perceptron_tfidf:.4f}"
288     )
289
290     print("SVM: LinearSVC")
291     print(f"{precision_svm_tfidf:.4f} {recall_svm_tfidf:.4f} {f1_svm_tfidf:.4f} {acc_svm_tfidf:.4f}")
292
293
294 if __name__ == "__main__":
295     main()
296
```

Overwriting HW1-CSCI544-wo-neg-sw.py
time: 50.9 ms (started: 2023-10-19 22:02:47 +00:00)

1 !python HW1-CSCI544-wo-neg-sw.py

Precision Recall F1-Score Accuracy
Perceptron
0.7637 0.8702 0.8135 0.7998
SVM: LinearSVC
0.8573 0.8602 0.8588 0.8581
time: 4min 1s (started: 2023-10-19 22:02:52 +00:00)

▼ Conclusion

Best Accuracies

Model	Accuracy	Features Used
Perceptron	0.8110	Word2Vec
LinearSVC	0.8321	Word2Vec
Perceptron	0.7998	TF-IDF

Model	Accuracy	Features Used
LinearSVC	0.8581	TF-IDF

1. LinearSVC outperforms Perceptron for both feature types (Word2Vec and TF-IDF).
- LinearSVC is better suited for this classification task compared to Perceptron.
2. When using Word2Vec features, both Perceptron and LinearSVC achieve lower accuracy compared to when using TF-IDF features.
- Word2Vec embeddings might not be as effective for this specific sentiment classification task as compared to TF-IDF vectors.
3. The LinearSVC model performs particularly well with TF-IDF features, achieving an accuracy of 85.81%.
- TF-IDF vectors are highly effective in capturing important information for sentiment classification in this dataset.

Overall, based on the provided performance metrics, it seems that TF-IDF features are more effective for this sentiment classification task compared to the Word2Vec embeddings. However, it's important to note that the effectiveness of features can vary depending on the specific dataset and task.

```
1 del balanced_df
2 del X_train, y_train, X_test, y_test

time: 372 µs (started: 2023-10-20 20:04:48 +00:00)
```

▼ Create Pytorch Dataset

- Custom pytorch dataset for on-the-fly processing and efficient resource utilization
- Each sample in this dataset includes embeddings and their corresponding target label. The label is adjusted by subtracting 1 from the label value in the DataFrame
- Using DataLoader's
 - Used to load and manage batches of data during the training process.
 - Handle tasks like shuffling, batching, and parallel data loading, making it easier to feed data to the model.

```
1 class AmazonReviewsSentimentDataset(Dataset):
2     def __init__(
3         self,
4         df, embeddings_col_name: str, label_col_name: str,
5         max_length=None, flatten: bool=True,
6         embedding_size: int=None, num_seq: int=None
7     ):
8         """Dataset class for Amazon Reviews Sentiment Analysis.
9
10        Args:
11            df (DataFrame): The input DataFrame containing the data.
12            embeddings_col_name (str): The column name for the embeddings.
13            label_col_name (str): The column name for the labels.
14            max_length (int, optional): Maximum length of embeddings (padding applied if needed).
15            flatten (bool, optional): Whether to flatten the embeddings or not.
16            embedding_size (int, optional): The size of each embedding.
17            num_seq (int, optional): The number of sequences (used when `flatten=False`).
18
19        Returns:
20            dict: A dictionary containing the embeddings and the target label.
21
22        Raises:
23            IndexError: If the index is out of bounds.
24        """
25
26        self.data = df
27        self.embeddings_col_name = embeddings_col_name
28        self.label_col_name = label_col_name
29        self.max_length = max_length
30        self.flatten = flatten
31        self.num_seq = num_seq
32        self.embedding_size = embedding_size
33
34    def __len__(self):
35        return len(self.data)
36
37    def __getitem__(self, idx):
38        if idx >= self.__len__():
39            raise IndexError
40
41        label = self.data.iloc[idx][self.label_col_name] - 1
42        embeddings = self.data.iloc[idx][self.embeddings_col_name]
43
44        # Pad embeddings to max_length if specified
45        if self.max_length is not None:
46            if len(embeddings) < self.max_length:
47                padding = np.zeros(self.max_length - len(embeddings), dtype=float)
48                embeddings = np.concatenate((embeddings, padding))
49
50        # Reshape embeddings if specified and flatten is False
51        if not self.flatten and self.num_seq is not None and self.embedding_size is not None:
52            embeddings = embeddings.reshape(self.num_seq, self.embedding_size)
53
54        return {
55            "embeddings": torch.tensor(embeddings, dtype=torch.float32),
56            "target": torch.tensor(label, dtype=torch.long)
57        }

time: 1.78 ms (started: 2023-10-20 22:55:50 +00:00)
```

```
1 TRAIN_BATCH_SIZE = 128
2 VALID_BATCH_SIZE = 64
3 TEST_BATCH_SIZE = 32
4 NUM_PARALLEL_WORKERS = multiprocessing.cpu_count()

time: 693 µs (started: 2023-10-20 21:38:30 +00:00)
```

▼ Training & Evaluation Functions

- `compute_accuracy` calculates the accuracy of model predictions given true labels.
- `train_loop_fn` handles one training epoch, updating the model's weights based on computed gradients.
- `eval_loop_fn` handles one validation epoch, computing the model's performance on the validation set.
- `train_and_evaluate` orchestrates the training process, saving checkpoints if specified. It reports metrics after each epoch. If a final model path is provided, it saves the model at the end.

```

1 def compute_accuracy(outputs, labels):
2     """
3     Computes the accuracy of the model's predictions.
4
5     Args:
6         outputs (torch.Tensor): The model's predictions.
7         labels (torch.Tensor): The true labels.
8
9     Returns:
10        float: The accuracy score.
11
12    """
13    predicted = torch.argmax(outputs.data, dim=1)
14
15    predicted = predicted.detach().cpu().numpy()
16    labels = labels.detach().cpu().numpy()
17
18    acc = accuracy_score(labels, predicted)
19    return acc
20
21 def train_loop_fn(data_loader, model, optimizer, loss_fn, device):
22     """
23     Performs one training epoch.
24
25     Args:
26         data_loader (DataLoader): The DataLoader for training data.
27         model (nn.Module): The neural network model.
28         optimizer (torch.optim): The optimizer for updating model weights.
29         loss_fn: The loss function.
30         device (torch.device): The device to perform computations.
31
32     Returns:
33        tuple: A tuple containing the training loss and accuracy.
34
35    """
36    model.train()
37    train_loss = 0.0
38    acc = []
39
40    for batch in tqdm(data_loader):
41        embeddings = batch['embeddings'].to(device, dtype=torch.float32, non_blocking=True)
42        labels = batch['target'].to(device, dtype=torch.long, non_blocking=True)
43
44        optimizer.zero_grad()
45
46        outputs = model(embeddings.float())
47        loss = loss_fn(outputs, labels)
48
49        loss.backward()
50        optimizer.step()
51
52        train_loss += loss.item()*len(labels)
53        acc.append(compute_accuracy(outputs, labels))
54
55    acc = sum(acc)/len(acc)
56    return train_loss, acc
57
58 def eval_loop_fn(data_loader, model, loss_fn, device):
59     """
60     Performs one evaluation epoch.
61
62     Args:
63         data_loader (DataLoader): The DataLoader for validation data.
64         model (nn.Module): The neural network model.
65         loss_fn: The loss function.
66         device (torch.device): The device to perform computations.
67
68     Returns:
69        tuple: A tuple containing the validation loss and accuracy.
70
71    """
72    valid_loss = 0.0
73    acc = []
74    model.eval()
75
76    for batch in data_loader:
77        embeddings = batch['embeddings'].to(device, dtype=torch.float32, non_blocking=True)
78        labels = batch['target'].to(device, dtype=torch.long, non_blocking=True)
79
80        outputs = model(embeddings.float())
81
82        loss = loss_fn(outputs, labels)
83        valid_loss += loss.item()*len(labels)
84
85        acc.append(compute_accuracy(outputs, labels))
86
87    acc = sum(acc)/len(acc)
88
89    return valid_loss, acc
90
91
92 def train_and_evaluate(
93     model,
94     train_data_loader, valid_data_loader,
95     optimizer, loss_fn,
96     device,
97     num_epochs,
98     checkpoint=False,
99     path="model.pt",
100     early_stopping_patience=5
101 ):
102     """
103     Trains and evaluates the model.
104
105     Args:
106         model (nn.Module): The neural network model.
107         train_data_loader (DataLoader): The DataLoader for training data.
108         valid_data_loader (DataLoader): The DataLoader for validation data.
109         optimizer (torch.optim): The optimizer for updating model weights.
110         loss_fn: The loss function.
111         device (torch.device): The device to perform computations.
112         num_epochs (int): The number of epochs.
113         checkpoint (bool, optional): Whether to save model checkpoints.
114         path (str, optional): The path to save the model.
115         early_stopping_patience (int, optional): Number of epochs to wait before early stopping.
116
117     Returns:
118         nn.Module: The best model.
119

```

```
120     """
121
122     # Create directory for saving checkpoint model states
123     if checkpoint:
124         dirname = path.split(".")[0]
125         checkpoint_path = os.path.join(dirname)
126         if os.path.exists(checkpoint_path):
127             shutil.rmtree(checkpoint_path)
128         os.makedirs(dirname)
129
130     best_loss = float('inf')
131     no_improvement_count = 0
132     best_model = None
133
134     for epoch in range(num_epochs):
135         # Train Step
136         train_loss, train_acc = train_loop_fn(
137             train_data_loader, model, optimizer, loss_fn, device
138         )
139
140         # Validation Step
141         valid_loss, valid_acc = eval_loop_fn(valid_data_loader, model, loss_fn, device)
142
143         train_loss /= len(train_data_loader.dataset)
144         valid_loss /= len(valid_data_loader.dataset)
145
146         epoch_log = (
147             f"Epoch {epoch+1}/{num_epochs},"
148             f" Train Accuracy={train_acc:.4f}, Validation Accuracy={valid_acc:.4f},"
149             f" Train Loss={train_loss:.4f}, Validation Loss={valid_loss:.4f}"
150         )
151         print(epoch_log)
152
153         # Check for improvement in validation loss
154         if valid_loss < best_loss:
155             # Save checkpoint if needed
156             if checkpoint:
157                 cp = os.path.join(checkpoint_path, f"{dirname}_epoch{epoch}_loss{valid_loss:.4f}.pt")
158                 torch.save(model.state_dict(), cp)
159                 print(f"Validation loss improved from {best_loss:.4f}-->{valid_loss:.4f}")
160                 print(f"Saved Checkpoint to '{cp}'")
161
162                 best_loss = valid_loss
163                 best_model = model
164                 no_improvement_count = 0
165             else:
166                 no_improvement_count += 1
167
168             # Early stopping condition
169             if no_improvement_count >= early_stopping_patience:
170                 print(f"No improvement for {early_stopping_patience} epochs. Stopping early.")
171                 break
172
173         if checkpoint:
174             # Save the best model
175             best_model_path = os.path.join(checkpoint_path, f"{dirname}-best.pt")
176             torch.save(best_model.state_dict(), best_model_path)
177             print(f"Saved best model to '{os.path.relpath(best_model_path)}'")
178
179         # Save current model
180         torch.save(model.state_dict(), path)
181
182     return best_model
183
184
185 def test_model(model, data_loader, device):
186     """
187     Tests the model on the test set.
188
189     Args:
190         model (nn.Module): The neural network model.
191         data_loader (DataLoader): The DataLoader for test data.
192         device (torch.device): The device to perform computations.
193
194     Returns:
195         tuple: A tuple containing the test accuracy and loss.
196
197     """
198     test_loss = 0.0
199     acc = []
200     loss_fn = nn.CrossEntropyLoss()
201
202     model.eval()
203     for batch in tqdm(data_loader):
204         embeddings = batch['embeddings'].to(device, dtype=torch.float32)
205         y_true = batch["target"].to(device, dtype=torch.long)
206
207         with torch.no_grad():
208             y_pred = model(embeddings)
209
210         loss = loss_fn(y_pred, y_true)
211         test_loss += loss.item()*len(y_true)
212
213         acc.append(compute_accuracy(y_pred, y_true))
214
215     acc = sum(acc)/len(acc)
216     test_loss = test_loss/len(data_loader.dataset)
217     return acc, test_loss
```

time: 2.22 ms (started: 2023-10-20 22:59:06 +00:00)

▼ Feedforward Neural Networks

```
1 class MLP(nn.Module):
2     def __init__(self, num_input_features, num_classes):
3         """
4         Multi-Layer Perceptron (MLP) for classification tasks.
5
6         Args:
7             num_input_features (int): Number of input features.
8             num_classes (int): Number of output classes.
9
10        """
11        super(MLP, self).__init__()
12        # Input size is 300 (Word2Vec dimensions)
13        self.fc1 = nn.Linear(num_input_features, 50)
```

```
14         self.fc2 = nn.Linear(50, 5)
15         # Output size is 2 for binary classification
16         self.fc3 = nn.Linear(5, num_classes)
17
18     def forward(self, x):
19         x = torch.relu(self.fc1(x))
20         x = torch.relu(self.fc2(x))
21         x = self.fc3(x)
22         return x
time: 860 µs (started: 2023-10-20 20:17:54 +00:00)
```

```
1 train_dataset = AmazonReviewsSentimentDataset(
2     train_df, embeddings_col_name="embeddings", label_col_name="sentiment"
3 )
4
5 valid_dataset = AmazonReviewsSentimentDataset(
6     valid_df, embeddings_col_name="embeddings", label_col_name="sentiment"
7 )
8
9 test_dataset = AmazonReviewsSentimentDataset(
10     test_df, embeddings_col_name="embeddings", label_col_name="sentiment"
11 )
12
13 train_data_loader = DataLoader(
14     train_dataset,
15     batch_size=TRAIN_BATCH_SIZE,
16     drop_last=True,
17     shuffle=True,
18     # num_workers=NUM_PARALLEL_WORKERS
19 )
20
21 valid_data_loader = DataLoader(
22     valid_dataset,
23     batch_size=VALID_BATCH_SIZE,
24     drop_last=False,
25     shuffle=False,
26     # num_workers=NUM_PARALLEL_WORKERS
27 )
28
29 test_data_loader = DataLoader(
30     test_dataset,
31     batch_size=TEST_BATCH_SIZE,
32     drop_last=False,
33     shuffle=False,
34     # num_workers=NUM_PARALLEL_WORKERS
35 )
time: 924 µs (started: 2023-10-20 21:38:49 +00:00)
```

▼ With average Word2Vec features

```
1 net = MLP(num_input_features=Word2VecConfig.MAX_LENGTH, num_classes=2).to(device)
2 criterion = nn.CrossEntropyLoss()
3 optimizer = optim.SGD(net.parameters(), lr=0.01)
4
5 model = train_and_evaluate(
6     model=net,
7     train_data_loader=train_data_loader,
8     valid_data_loader=valid_data_loader,
9     optimizer=optimizer,
10    loss_fn=criterion,
11    device=device,
12    num_epochs=25,
13    checkpoint=True,
14    path="mlp_w_avg_w2v_feat_v3.pt"
15 )
```

100%

624/624 [00:18<00:00, 38.44it/s]

Epoch 1/25, Train Accuracy=0.5131, Validation Accuracy=0.5611, Train Loss=0.6950, Validation Loss=0.6919
Validation loss improved from inf--->0.6919
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch0_loss0.6919.pt'

100%

624/624 [00:16<00:00, 38.91it/s]

Epoch 2/25, Train Accuracy=0.6316, Validation Accuracy=0.6768, Train Loss=0.6911, Validation Loss=0.6905
Validation loss improved from 0.6919--->0.6905
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch1_loss0.6905.pt'

100%

624/624 [00:17<00:00, 37.97it/s]

Epoch 3/25, Train Accuracy=0.6516, Validation Accuracy=0.6777, Train Loss=0.6894, Validation Loss=0.6884
Validation loss improved from 0.6905--->0.6884
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch2_loss0.6884.pt'

100%

624/624 [00:17<00:00, 30.38it/s]

Epoch 4/25, Train Accuracy=0.6742, Validation Accuracy=0.6149, Train Loss=0.6868, Validation Loss=0.6851
Validation loss improved from 0.6884--->0.6851
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch3_loss0.6851.pt'

100%

624/624 [00:16<00:00, 38.13it/s]

Epoch 5/25, Train Accuracy=0.6457, Validation Accuracy=0.6580, Train Loss=0.6818, Validation Loss=0.6777
Validation loss improved from 0.6851--->0.6777
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch4_loss0.6777.pt'

100%

624/624 [00:18<00:00, 36.78it/s]

Epoch 6/25, Train Accuracy=0.6730, Validation Accuracy=0.6678, Train Loss=0.6723, Validation Loss=0.6655
Validation loss improved from 0.6777--->0.6655
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch5_loss0.6655.pt'

100%

624/624 [00:17<00:00, 38.92it/s]

Epoch 7/25, Train Accuracy=0.6914, Validation Accuracy=0.6881, Train Loss=0.6564, Validation Loss=0.6448
Validation loss improved from 0.6655--->0.6448
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch6_loss0.6448.pt'

100%

624/624 [00:17<00:00, 38.77it/s]

Epoch 8/25, Train Accuracy=0.7160, Validation Accuracy=0.7258, Train Loss=0.6308, Validation Loss=0.6127
Validation loss improved from 0.6448--->0.6127
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch7_loss0.6127.pt'

100%

624/624 [00:17<00:00, 38.84it/s]

Epoch 9/25, Train Accuracy=0.7429, Validation Accuracy=0.7512, Train Loss=0.5944, Validation Loss=0.5711
Validation loss improved from 0.6127--->0.5711
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch8_loss0.5711.pt'

100%

624/624 [00:17<00:00, 39.25it/s]

Epoch 10/25, Train Accuracy=0.7628, Validation Accuracy=0.7694, Train Loss=0.5528, Validation Loss=0.5292
Validation loss improved from 0.5711--->0.5292
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch9_loss0.5292.pt'

100%

624/624 [00:17<00:00, 29.82it/s]

Epoch 11/25, Train Accuracy=0.7751, Validation Accuracy=0.7859, Train Loss=0.5159, Validation Loss=0.4955
Validation loss improved from 0.5292--->0.4955
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch10_loss0.4955.pt'

100%

624/624 [00:16<00:00, 37.17it/s]

Epoch 12/25, Train Accuracy=0.7846, Validation Accuracy=0.7918, Train Loss=0.4887, Validation Loss=0.4720
Validation loss improved from 0.4955--->0.4720
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch11_loss0.4720.pt'

100%

624/624 [00:17<00:00, 37.59it/s]

Epoch 13/25, Train Accuracy=0.7917, Validation Accuracy=0.7959, Train Loss=0.4699, Validation Loss=0.4560
Validation loss improved from 0.4720--->0.4560
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch12_loss0.4560.pt'

100%

624/624 [00:17<00:00, 37.73it/s]

Epoch 14/25, Train Accuracy=0.7967, Validation Accuracy=0.8017, Train Loss=0.4565, Validation Loss=0.4443
Validation loss improved from 0.4560--->0.4443
Saving Checkpoint to 'mlp_w_avg_w2v_feat_v3/mlp_w_avg_w2v_feat_v3_epoch13_loss0.4443.pt'

100%

624/624 [00:17<00:00, 38.39it/s]

Overall Accuracy on Test Set

Saving Checkpoint to 'mln w avg w2v feat v3/mln w avg w2v feat v3_epoch14_loss0.4353.pt'

1 path_to_saved_model = 'mlp_w_avg_w2v_feat_v3.pt'
2 model = MLP(num_input_features=Word2VecConfig.MAX_LENGTH, num_classes=2)
3 model.load_state_dict(torch.load(path_to_saved_model, map_location=device))
4
5 acc, loss = test_model(model, test_data_loader, device)
6 print("Accuracy (Test Dataset):", round(acc,4))

100%

157/157 [00:02<00:00, 95.96it/s]

Accuracy (Test Dataset): 0.8298
time: 2.45 s (started: 2023-10-20 21:43:29 +00:00)
Validation loss improved from 0.4370 --> 0.4176

With top 10 Word2Vec features

- Embeddings are padded for maintaining consistent input dimensions across different samples in a batch.

Saving Checkpoint to 'mlp w avg w2v feat v3/mlp w avg w2v feat v3_epoch18_loss0.4137.pt'

1 train_dataset = AmazonReviewsSentimentDataset(
2 train_df,
3 embeddings_col_name="embeddings_top_10",
4 label_col_name="sentiment",
5 max_length=3000,
6 flatten=True
7)
8
9 valid_dataset = AmazonReviewsSentimentDataset(
10 valid_df,
11 embeddings_col_name="embeddings_top_10",
12 label_col_name="sentiment",
13 max_length=3000,
14 flatten=True
15)
16
17 test_dataset = AmazonReviewsSentimentDataset(
18 test_df,
19 embeddings_col_name="embeddings_top_10",
20 label_col_name="sentiment",
21 max_length=3000,
22 flatten=True
23)
24
25 train_data_loader = DataLoader(
26 train_dataset,
27 batch_size=TRAIN_BATCH_SIZE,
28 drop_last=True,
29 shuffle=True,
30)
31
32 valid_data_loader = DataLoader(
33 valid_dataset,
34 batch_size=VALID_BATCH_SIZE,

```
35     drop_last=False,
36     shuffle=False,
37 )
38
39 test_data_loader = DataLoader(
40     test_dataset,
41     batch_size=TEST_BATCH_SIZE,
42     drop_last=False,
43     shuffle=False,
44 )
```

time: 1.29 ms (started: 2023-10-20 21:46:47 +00:00)

```
1 net2 = MLP(num_input_features=3000, num_classes=2).to(device)
2 criterion = nn.CrossEntropyLoss()
3 optimizer = optim.SGD(net2.parameters(), lr=0.01)
4
5
6 model2 = train_and_evaluate(
7     model=net2,
8     train_data_loader=train_data_loader,
9     valid_data_loader=valid_data_loader,
10    optimizer=optimizer,
11    loss_fn=criterion,
12    device=device,
13    num_epochs=30,
14    checkpoint=True,
15    path="mlp_w_top10_w2v_feat_v2.pt"
16 )
```

```
100%                               624/624 [00:18<00:00, 36.25it/s]

Epoch 1/25, Train Accuracy=0.5929, Validation Accuracy=0.6606, Train Loss=0.6864, Validation Loss=0.6760
Validation loss improved from inf-->0.6760
Saved Checkpoint to 'mlp_w_top10_w2v_feat_v2/mlp_w_top10_w2v_feat_v2_epoch0_loss0.6760.pt'

100%                               624/624 [00:18<00:00, 33.14it/s]

Epoch 2/25, Train Accuracy=0.6798, Validation Accuracy=0.7102, Train Loss=0.6556, Validation Loss=0.6253
Validation loss improved from 0.6760-->0.6253
Saved Checkpoint to 'mlp_w_top10_w2v_feat_v2/mlp_w_top10_w2v_feat_v2_epoch1_loss0.6253.pt'

100%                               624/624 [00:18<00:00, 34.66it/s]

Epoch 3/25, Train Accuracy=0.7199, Validation Accuracy=0.7278, Train Loss=0.5863, Validation Loss=0.5492
Validation loss improved from 0.6253-->0.5492
Saved Checkpoint to 'mlp_w_top10_w2v_feat_v2/mlp_w_top10_w2v_feat_v2_epoch2_loss0.5492.pt'

100%                               624/624 [00:18<00:00, 36.23it/s]
```

Overall Accracy on Test Set

```
Saved Checkpoint to 'mlp w top10 w2v feat v2/mlp w top10 w2v feat v2 epoch3 loss0.5210.pt'

1 path_to_saved_model = 'mlp_w_top10_w2v_feat_v2.pt'
2 model = MLP(num_input_features=3000, num_classes=2)
3 model.load_state_dict(torch.load(path_to_saved_model, map_location=device))
4
5 acc, loss = test_model(model, test_data_loader, device)
6 print("Accuracy (Test Dataset):", round(acc,4))
```

```
100%                               157/157 [00:03<00:00, 58.79it/s]

Accuracy (Test Dataset): 0.7761
time: 3.41 s (started: 2023-10-20 21:56:52 +00:00)

Validation loss improved from 0.5041-->0.5000
```

Comparision with Simple Model

The LinearSVC model trained on TF-IDF features was the most effective in this scenario, outperforming both simple models and MLP models trained with Word2Vec embeddings.

Conclusion

1. Feature Importance:

- The choice of features significantly impacts model performance.
- In this case, TF-IDF features proved to be the most informative for sentiment analysis, as evidenced by the high accuracy achieved by LinearSVC with TF-IDF.

2. Complexity vs. Performance:

- Simple models like Perceptron and LinearSVC can sometimes outperform more complex models.
- This is evident in the case where LinearSVC with TF-IDF outperformed the MLP models.

3. Embedding Selection:

- Not all embeddings are equally effective. The choice of Word2Vec embeddings, particularly using the average vectors, yielded competitive results, showcasing the importance of using quality word embeddings.

4. Dimensionality Matters:

- Using only the top 10 Word2Vec embeddings didn't capture enough information for sentiment analysis.
- It's important to consider the dimensionality of the embeddings and how well they represent the underlying semantics.

```
Epoch 14/25, Train Accuracy=0.7784, Validation Accuracy=0.7602, Train Loss=0.4664, Validation Loss=0.4965
```

▼ Recurrent Neural Networks

```
1 class RNNModel(nn.Module):
2     def __init__(
3         self, input_size, hidden_size, num_layers, output_size, model_type="rnn"
4     ):
5         """
6         Recurrent Neural Network (RNN) model for sequence data processing.
7
8         Args:
9             input_size (int): Dimension of the input features.
10            hidden_size (int): Number of units in the hidden layers.
11            num_layers (int): Number of recurrent layers.
12            output_size (int): Number of output classes.
13            model_type (str, optional): Type of RNN ('rnn', 'gru', or 'lstm'). Defaults to 'rnn'.
14
15        """
16        super(RNNModel, self).__init__()
17
18        self.hidden_size = hidden_size
19        self.num_layers = num_layers
20        self.model_type = model_type
21
22        # Initialize the recurrent layer based on model_type
23        if model_type == "gru":
24            self.layer = nn.GRU(input_size, hidden_size, num_layers, batch_first=True, dropout=0.3)
25        elif model_type == "lstm":
26            self.layer = nn.LSTM(input_size, hidden_size, num_layers, batch_first=True, dropout=0.3)
27        else:
28            self.layer = nn.RNN(input_size, hidden_size, num_layers, batch_first=True, dropout=0.3)
29
30        # dropout layer to prevent overfitting
31        self.dropout = nn.Dropout(0.3)
32
33        # Fully connected layer for final prediction
34        self.fc = nn.Linear(hidden_size, output_size)
35
36    def forward(self, x):
37        batch_size = x.size(0)
38        hidden = self.init_hidden(batch_size)
39
40        # Pass input through the recurrent layer
41        out, _ = self.layer(x, hidden)
42
43        # Stack up the model output
44        # out = out.contiguous().view(-1, self.hidden_size)
45
46        # Use the output from the last time step
47        out = out[:, -1, :]
48
49        # out = self.dropout(out)
50
51        # Apply fully connected layer for final prediction
52        out = self.fc(out)
53        return out
54
```



```
55     def init_hidden(self, batch_size):
56         if self.model_type == "lstm":
57             hidden = (
58                 torch.zeros(self.num_layers, batch_size, self.hidden_size).to(device),
59                 torch.zeros(self.num_layers, batch_size, self.hidden_size).to(device)
60             )
61         else:
62             hidden = torch.zeros(self.num_layers, batch_size, self.hidden_size).to(device)
63         return hidden
```

time: 1.65 ms (started: 2023-10-20 21:57:09 +00:00)

```
1 train_dataset = AmazonReviewsSentimentDataset(
2     train_df,
3     embeddings_col_name="embeddings_top_10",
4     label_col_name="sentiment",
5     max_length=3000,
6     flatten=False,
7     embedding_size=Word2VecConfig.EMBEDDING_SIZE,
8     num_seq=10
9 )
10
11 valid_dataset = AmazonReviewsSentimentDataset(
12     valid_df,
13     embeddings_col_name="embeddings_top_10",
14     label_col_name="sentiment",
15     max_length=3000,
16     flatten=False,
17     embedding_size=Word2VecConfig.EMBEDDING_SIZE,
18     num_seq=10
19 )
20
21 test_dataset = AmazonReviewsSentimentDataset(
22     test_df,
23     embeddings_col_name="embeddings_top_10",
24     label_col_name="sentiment",
25     max_length=3000,
26     flatten=False,
27     embedding_size=Word2VecConfig.EMBEDDING_SIZE,
28     num_seq=10
29 )
30
31 train_data_loader = DataLoader(
32     train_dataset,
33     batch_size=TRAIN_BATCH_SIZE,
34     drop_last=True,
35     shuffle=True,
36 )
37
38 valid_data_loader = DataLoader(
39     valid_dataset,
40     batch_size=VALID_BATCH_SIZE,
41     drop_last=False,
42     shuffle=False,
43 )
44
45 test_data_loader = DataLoader(
46     test_dataset,
47     batch_size=TEST_BATCH_SIZE,
48     drop_last=False,
49     shuffle=False,
50 )
```

time: 1.11 ms (started: 2023-10-20 23:05:20 +00:00)

```
1 input_size = 300
2 hidden_size = 10
3 output_size = 2
4 num_layers = 10
```

time: 473 µs (started: 2023-10-20 23:05:22 +00:00)

▼ Simple RNN

```
1 net3 = RNNModel(input_size, hidden_size, num_layers, output_size, model_type="rnn").to(device)
2 criterion = nn.CrossEntropyLoss()
3 optimizer = torch.optim.Adam(net3.parameters(), lr=0.001)
4
5 model3 = train_and_evaluate(
6     model=net3,
7     train_data_loader=train_data_loader,
8     valid_data_loader=valid_data_loader,
9     optimizer=optimizer,
10    loss_fn=criterion,
11    device=device,
12    num_epochs=20,
13    checkpoint=True,
14    path="simple_rnn_w2v_feat_v2.pt"
15 )
```

```
100%                                624/624 [00:35<00:00, 26.98it/s]
Epoch 1/20, Train Accuracy=0.6423, Validation Accuracy=0.7367, Train Loss=0.6211, Validation Loss=0.5484
Validation loss improved from inf--->0.5484
Saved Checkpoint to 'simple_rnn_w2v_feat_v2/simple_rnn_w2v_feat_v2_epoch0_loss0.5484.pt'
100%                                624/624 [00:28<00:00, 20.17it/s]
Epoch 2/20, Train Accuracy=0.7385, Validation Accuracy=0.7457, Train Loss=0.5413, Validation Loss=0.5386
Validation loss improved from 0.5484--->0.5386
Saved Checkpoint to 'simple_rnn_w2v_feat_v2/simple_rnn_w2v_feat_v2_epoch1_loss0.5386.pt'
100%                                624/624 [00:26<00:00, 26.29it/s]
Epoch 3/20, Train Accuracy=0.7492, Validation Accuracy=0.7528, Train Loss=0.5260, Validation Loss=0.5217
Validation loss improved from 0.5386--->0.5217
Saved Checkpoint to 'simple_rnn_w2v_feat_v2/simple_rnn_w2v_feat_v2_epoch2_loss0.5217.pt'
100%                                624/624 [00:37<00:00, 25.14it/s]
Epoch 4/20, Train Accuracy=0.7559, Validation Accuracy=0.7570, Train Loss=0.5182, Validation Loss=0.5154
Validation loss improved from 0.5217--->0.5154
Saved Checkpoint to 'simple_rnn_w2v_feat_v2/simple_rnn_w2v_feat_v2_epoch3_loss0.5154.pt'
100%                                624/624 [00:29<00:00, 20.49it/s]
Epoch 5/20, Train Accuracy=0.7605, Validation Accuracy=0.7617, Train Loss=0.5119, Validation Loss=0.5129
Validation loss improved from 0.5154--->0.5129
Saved Checkpoint to 'simple_rnn_w2v_feat_v2/simple_rnn_w2v_feat_v2_epoch4_loss0.5129.pt'
100%                                624/624 [00:26<00:00, 25.89it/s]
Epoch 6/20, Train Accuracy=0.7609, Validation Accuracy=0.7610, Train Loss=0.5079, Validation Loss=0.5058
Validation loss improved from 0.5129--->0.5058
Saved Checkpoint to 'simple_rnn_w2v_feat_v2/simple_rnn_w2v_feat_v2_epoch5_loss0.5058.pt'
100%                                624/624 [00:26<00:00, 25.14it/s]
Epoch 7/20, Train Accuracy=0.7624, Validation Accuracy=0.7595, Train Loss=0.5046, Validation Loss=0.5097
100%                                624/624 [00:26<00:00, 18.12it/s]
Epoch 8/20, Train Accuracy=0.7641, Validation Accuracy=0.7600, Train Loss=0.5016, Validation Loss=0.4996
```

Overall Accuracy on Test Set

```
1 path_to_saved_model = 'simple_rnn_w2v_feat_v2.pt'
2 model = RNNModel(input_size, hidden_size, num_layers, output_size, model_type="rnn")
3 model.load_state_dict(torch.load(path_to_saved_model, map_location=device))
4
5 acc, loss = test_model(model, test_data_loader, device)
6 print("Accuracy (Test Dataset):", round(acc,4))
```

```
100%                                157/157 [00:01<00:00, 113.81it/s]
Accuracy (Test Dataset): 0.7765
time: 1.49 s (started: 2023-10-20 23:23:30 +00:00)
```

▼ GRU

```
1 net4 = RNNModel(input_size, hidden_size, num_layers, output_size, model_type="gru").to(device)
2 criterion = nn.CrossEntropyLoss()
3 optimizer = torch.optim.Adam(net4.parameters(), lr=0.001)
4
5 model4 = train_and_evaluate(
6     model=net4,
7     train_data_loader=train_data_loader,
8     valid_data_loader=valid_data_loader,
9     optimizer=optimizer,
10    loss_fn=criterion,
11    device=device,
12    num_epochs=20,
13    checkpoint=True,
14    path="gru_w2v_feat_v2.pt"
15 )
```

100%

624/624 [00:42<00:00, 16.88it/s]

Epoch 1/20, Train Accuracy=0.6928, Validation Accuracy=0.7528, Train Loss=0.5734, Validation Loss=0.5106
Validation loss improved from inf--->0.5106
Saved Checkpoint to 'gru_w2v_feat_v2/gru_w2v_feat_v2_epoch0_loss0.5106.pt'

100%

624/624 [00:37<00:00, 15.09it/s]

Epoch 2/20, Train Accuracy=0.7611, Validation Accuracy=0.7681, Train Loss=0.5031, Validation Loss=0.4882
Validation loss improved from 0.5106--->0.4882
Saved Checkpoint to 'gru_w2v_feat_v2/gru_w2v_feat_v2_epoch1_loss0.4882.pt'

100%

624/624 [00:38<00:00, 16.97it/s]

Epoch 3/20, Train Accuracy=0.7744, Validation Accuracy=0.7754, Train Loss=0.4803, Validation Loss=0.4765
Validation loss improved from 0.4882--->0.4765
Saved Checkpoint to 'gru_w2v_feat_v2/gru_w2v_feat_v2_epoch2_loss0.4765.pt'

100%

624/624 [00:38<00:00, 17.81it/s]

Epoch 4/20, Train Accuracy=0.7806, Validation Accuracy=0.7785, Train Loss=0.4677, Validation Loss=0.4810

100%

624/624 [00:38<00:00, 13.10it/s]

Epoch 5/20, Train Accuracy=0.7858, Validation Accuracy=0.7834, Train Loss=0.4601, Validation Loss=0.4533
Validation loss improved from 0.4765--->0.4533
Saved Checkpoint to 'gru_w2v_feat_v2/gru_w2v_feat_v2_epoch4_loss0.4533.pt'

100%

624/624 [00:38<00:00, 17.16it/s]

Epoch 6/20, Train Accuracy=0.7910, Validation Accuracy=0.7794, Train Loss=0.4504, Validation Loss=0.4791

100%

624/624 [00:38<00:00, 17.37it/s]

Overall Accuracy on Test Set

```
1 path_to_saved_model = 'gru_w2v_feat_v2.pt'
2 model = RNNModel(input_size, hidden_size, num_layers, output_size, model_type="gru")
3 model.load_state_dict(torch.load(path_to_saved_model, map_location=device))
4
5 acc, loss = test_model(model, test_data_loader, device)
6 print("Accuracy (Test Dataset):", round(acc,4))
```

100%

157/157 [00:02<00:00, 61.83it/s]

Accuracy (Test Dataset): 0.8063
time: 2.32 s (started: 2023-10-20 22:36:42 +00:00)

▼ LSTM

100%

624/624 [00:38<00:00, 16.88it/s]

```
1 net5 = RNNModel(input_size, hidden_size, num_layers, output_size, model_type="lstm").to(device)
2 criterion = nn.CrossEntropyLoss()
3 optimizer = torch.optim.Adam(net5.parameters(), lr=0.001)
4
5 model5 = train_and_evaluate(
6     model=net5,
7     train_data_loader=train_data_loader,
8     valid_data_loader=valid_data_loader,
9     optimizer=optimizer,
10    loss_fn=criterion,
11    device=device,
12    num_epochs=20,
13    checkpoint=True,
14    path="lstm_w2v_feat_v2.pt"
15 )
```

100% 624/624 [00:31<00:00, 23.71it/s]

Epoch 1/20, Train Accuracy=0.6529, Validation Accuracy=0.7387, Train Loss=0.6076, Validation Loss=0.5329
Validation loss improved from inf-->0.5329
Saved Checkpoint to 'lstm_w2v_feat_v2/lstm_w2v_feat_v2_epoch0_loss0.5329.pt'

100% 624/624 [00:27<00:00, 23.73it/s]

Epoch 2/20, Train Accuracy=0.7565, Validation Accuracy=0.7735, Train Loss=0.5129, Validation Loss=0.4886
Validation loss improved from 0.5329-->0.4886
Saved Checkpoint to 'lstm_w2v_feat_v2/lstm_w2v_feat_v2_epoch1_loss0.4886.pt'

100% 624/624 [00:28<00:00, 23.71it/s]

Epoch 3/20, Train Accuracy=0.7732, Validation Accuracy=0.7785, Train Loss=0.4826, Validation Loss=0.4754
Validation loss improved from 0.4886-->0.4754
Saved Checkpoint to 'lstm_w2v_feat_v2/lstm_w2v_feat_v2_epoch2_loss0.4754.pt'

100% 624/624 [00:29<00:00, 23.78it/s]

Epoch 4/20. Train Accuracv=0.7800. Validation Accuracv=0.7858. Train Loss=0.4669. Validation Loss=0.4692

Overall Accuracy on Test Set

```
100% 624/624 [00:27<00:00, 23.80it/s]
1 path_to_saved_model = 'lstm_w2v_feat_v2.pt'
2 model = RNNModel(input_size, hidden_size, num_layers, output_size, model_type="lstm")
3 model.load_state_dict(torch.load(path_to_saved_model, map_location=device))
4
5 acc, loss = test_model(model, test_data_loader, device)
6 print("Accuracy (Test Dataset):", round(acc,4))

100% 157/157 [00:01<00:00, 101.68it/s]

Accuracy (Test Dataset): 0.8023
time: 1.79 s (started: 2023-10-20 22:52:48 +00:00)

Saved Checkpoint to 'lstm_w2v_feat_v2/lstm_w2v_feat_v2_epoch6_loss0.4518.pt'
```

Conclusion

- 1. **Feature Representations:**
 - TF-IDF outperforms Word2Vec across all models.
 - Averaged Word2Vec is better than Concatenated Word2Vec.
- 2. **Model Comparisons:**
 - SVM outperforms Perceptron consistently.
 - MLP with averaged Word2Vec performs better than RNN, GRU, and LSTM with Word2Vec.
- 3. **Recurrent Models:**
 - RNN, GRU, and LSTM show similar performance with Word2Vec embeddings.
- 4. **Overall Performance:**
 - Highest accuracy (~87%) is achieved with SVM using TF-IDF.

Other:

- Averaging Word2Vec embeddings seems a more effective representation
- SVM model is better at capturing the non-linear relationships in the data compared to the Perceptron
- TF-IDF may capture important information more effectively than Word2Vec embeddings

Report

- Results from one of the runs

Model	Accuracy	Features
SVM(LinearSVC)	0.8581	TF-IDF
SVM(LinearSVC)	0.8321	mean word2vec
MLP	0.8298	mean word2vec
Perceptron	0.8110	mean word2vec
GRU	0.8063	top 10 word2vec
LSTM	0.8023	top 10 word2vec
Perceptron	0.7998	TF-IDF
Simple RNN	0.7765	top 10 word2vec
MLP	0.7761	top 10 word2vec

References

- 1. <https://www.kaggle.com/code/abhishek/bert-multi-lingual-tpu-training-8-cores>
- 2. <https://www.kaggle.com/mishra1993/pytorch-multi-layer-perceptron-mnist>
- 3. <https://www.kaggle.com/code/arunmohan003/sentiment-analysis-using-lstm-pytorch>
- 4. <https://pytorch.org/docs/stable/generated/torch.nn.RNN.html>
- 5. https://pytorch.org/tutorials/intermediate/char_rnn_classification_tutorial.html
- 6. https://radimrehurek.com/gensim/auto_examples/tutorials/run_word2vec.html

THE END