

# Assignment 4: Sentiment Analysis using Feedforward Neural Networks

Ellie Larence, Matthew Lloyd Macias, Jamerson Tenorio, Leo Wang

April 25, 2025

## 2 Base FNN Model

The base FNN model constructed by Matthew Macias consists of 3 hidden layers, where the first two hidden layers have half the number of the initial input layer, and the final hidden layer with a fourth of the initial input layer. The FNN applies a linear transformation, on each layer, using ReLU as the activation function. The dataset provided to the FNN was first split into a 70-30 training-testing data split, and then vectorized using tfidf vectorization utilizing the preprocessing and Porter tokenization functions shown in the textbook and during class on the slides. The testing and training features were then initially normalized before being passed into a tensor dataset that would afterwards be utilized by the FNN. The labels were untouched, but split in the same manner as the features. The optimizer and loss function our FNN uses is the SGD optimizer, and the Cross Entropy Loss function. No L2 Regularization parameter or learning rate was specified at this time, as to be able to tune these in Task 3. Our Base FNN also only trains over 5 epochs.

The split was modified by Leo Wang to ensure correct splitting of training and testing data. Leo Wang also additionally added the total and average loss calculations, as well as time measurement to the training loop. This Base FNN is found in the `generate_data.py` script provided in the Assignment zip file, and outlined here in the following code snippet.

```
1 import pyprind
2 import pandas as pd
3 import numpy as np
4 import os
5 import sys
6 import re
7 import time
8 from sklearn.metrics import accuracy_score
9 from sklearn.feature_extraction.text import TfidfVectorizer
10 from sklearn.model_selection import train_test_split
11
12 # Don't think these will be used for me (Matt)
13 # But maybe for others
14 from sklearn.model_selection import GridSearchCV
15 from sklearn.pipeline import Pipeline
16 from sklearn.linear_model import LogisticRegression
17
18 import torch
19 import torch.nn as nn
20 from torch.utils.data import DataLoader
21 from torch.utils.data import TensorDataset
22
23 import nltk
24 nltk.download('stopwords')
25
26 from nltk.corpus import stopwords
27 from nltk.stem.porter import PorterStemmer
28
29 # Global seeds
30 torch.manual_seed(42)
31 np.random.seed(42)
32
33 # Global variables
34 csv_file_name = './movie_data.csv'
35 stop = stopwords.words('english')
36 porter = PorterStemmer()
37
```

```

38 # Generates a dataframe object from the movie reviews.
39 # returns
40 # - df - A dataframe object containing all reviews and their sentiments
41 def generate_df_from_reviews():
42     # If the csv file has not been generated yet,
43     # generate from the aclImdb directory.
44     if not (os.path.exists(csv_file_name)):
45         current_dir = os.path.dirname(__file__)
46
47         data_dir_name = "aclImdb"
48         dataset_dir = os.path.join(current_dir, data_dir_name)
49         labels = {'pos':1, 'neg':0}
50         # pbar can be commented out
51         pbar = pyprind.ProgBar(50000, stream=sys.stdout)
52         df = pd.DataFrame()
53         for s in ('test', 'train'):
54             for l in ('pos', 'neg'):
55                 path = os.path.join(dataset_dir, s, l)
56                 for file in sorted(os.listdir(path)):
57                     with open(os.path.join(path, file), 'r', encoding='utf-8') as infile:
58                         txt = infile.read()
59                         df_extension = pd.DataFrame([[txt, labels[l]]])
60                         df = pd.concat([df, df_extension], ignore_index=True)
61                         # pbar can be commented out
62                         pbar.update()
63
64         # Label columns
65         df.columns = ['review', 'sentiment']
66         # Shuffle
67         df = df.reindex(np.random.permutation(df.index))
68         # Store to csv
69         df.to_csv(csv_file_name, index=False, encoding='utf-8')
70     # Else, data has already been written, so read in from existing csv
71     else:
72         df = pd.read_csv(csv_file_name, encoding='utf-8')
73         df = df.rename(columns={"0": "review", "1": "sentiment"})
74
75     # Return resulting dataframe. Generated from seed 42
76     return df
77
78
79 # Preprocessor to clean the data (From Textbook/slides)
80 def preprocessor(text):
81     text = re.sub('<[^>]*>', '', text)
82     # emoticons = re.findall('(?:[:|;|=)(?:-)?(?:\)|\(|D|P)', text)
83     # text = (re.sub('[\W]+', ' ', text.lower()) + ' ').join(emoticons).replace('-', '')
84     # adding r for raw string to circumvent SyntaxWarning
85     emoticons = re.findall(r'(?:[:|;|=)(?:-)?(?:\)|\(|D|P)', text)
86     text = (re.sub(r'[\W]+', ' ', text.lower()) + ' ').join(emoticons).replace('-', '')
87     return text
88
89
90 # Tokenizer (From Textbook/slides)
91 def tokenizer_porter(text):
92     return [porter.stem(word) for word in text.split()]

```

```

93
94
95 # Process the dataframe generated from the reviews
96 def process_data():
97     start_time = time.time()
98
99     # Generate DataFrame Object where data is stored
100     print("\n1. Loading dataset...")
101     df = generate_df_from_reviews()
102     print(f"Total number of reviews in dataset: {len(df)}")
103
104     # 70-30 split
105     # X_train = df.loc[:35000, 'review'].values
106     # X_test = df.loc[15000:, 'review'].values
107     # y_train = df.loc[:35000, 'sentiment'].values
108     # y_test = df.loc[15000:, 'sentiment'].values
109     # ^ Old, bad split, caused very bad accuracy ^
110     # data leakage due to overlap from row 15000 through 35000 between training and
111         testing data!
112
113     # New, proper 70-30 split from Leo Wang
114     split_idx = 35000
115     X_train = df.loc[:split_idx-1, 'review'].values
116     X_test = df.loc[split_idx:, 'review'].values
117     y_train = df.loc[:split_idx-1, 'sentiment'].values
118     y_test = df.loc[split_idx:, 'sentiment'].values
119
120     print(f"Training set size: {len(X_train)}")
121     print(f"Test set size: {len(X_test)}")
122
123     # Limit max features to roughly half of what it
124     # would originally generate with this vectorizer
125     # Change max_features to include more features for
126     # processing. 50000 features went to about 2-3 hours
127     # Original size from a very initial test was around
128     # 100000 features. If you change this, change the scale
129     # of the FNN and change the input in net.forward()
130     # to also match this number.
131     tfidf = TfidfVectorizer(max_features=10000, strip_accents=None, lowercase=False,
132         preprocessor=preprocessor, tokenizer=tokenizer_porter, stop_words='english')
133
134     # Vectorize the review data and turn it into something
135     # that the TensorDataset can ingest
136     # Also vectorize test data to test on finished model
137     tfidf_reviews = tfidf.fit_transform(X_train).toarray()
138
139     # tfidf_testing = tfidf.fit_transform(X_test).toarray()
140     tfidf_testing = tfidf.transform(X_test).toarray()
141     # doesn't need to fit the vectorizer again, test data should be vectorized using the
142     # same vocab and transform learned from the training data
143
144     tfidf_reviews_norm = (tfidf_reviews - np.mean(tfidf_reviews)) / np.std(tfidf_reviews)
145     tfidf_testing_norm = (tfidf_testing - np.mean(tfidf_testing)) / np.std(tfidf_testing)
146
147     # Float for the reviews, as the array is likely to have floats rather than ints

```

```

145     # long for the sentiments to properly gauge positive or negative
146     tfidf_reviews_norm = torch.from_numpy(tfidf_reviews_norm).float()
147     tfidf_testing_norm = torch.from_numpy(tfidf_testing_norm).float()
148     y_train = torch.from_numpy(y_train).long()
149
150
151     # Initialize the dataset and load it into the dataloader
152     train_ds = TensorDataset(tfidf_reviews_norm, y_train)
153     review_data = DataLoader(train_ds, batch_size=100, shuffle=True)
154
155
156     # initialize FNN. Model based off of Linear ReLU from class
157     print("\n4. Initializing neural network...")
158     net = torch.nn.Sequential(
159         torch.nn.Linear(10000, 5000),
160         torch.nn.ReLU(),
161         torch.nn.Linear(5000, 5000),
162         torch.nn.ReLU(),
163         torch.nn.Linear(5000, 5000),
164         torch.nn.ReLU(),
165         torch.nn.Linear(5000, 2500),
166         torch.nn.ReLU(),
167         torch.nn.Linear(2500, 2),
168         torch.nn.Softmax(dim=1))
169
170     optimizer = torch.optim.SGD(net.parameters(), lr=0.1)
171     L = torch.nn.CrossEntropyLoss()
172
173     # Train. It's doing something and it takes a while. I got to take a shower and it
174     # barely finished the
175     # Second epoch. So to finish 5 epochs, likely 1 hour and 30 minutes.
176     print("\n5. Starting training...")
177     print("Epoch_|_Training_Loss_|_Test_Accuracy_|_Time")
178     print("-" * 40)
179
180     for epoch in range(5):
181         epoch_start = time.time()
182         net.train()
183         total_loss = 0
184         batch_count = 0
185
186         for (x, y) in review_data:
187             output = net.forward(x.view(-1,10000))
188             loss = L(output, y)
189             total_loss += loss.item()
190             batch_count += 1
191             loss.backward()
192             optimizer.step()
193             net.zero_grad()
194
195         avg_loss = total_loss / batch_count
196
197         net.eval()
198         with torch.no_grad():
199             y_pred = net(tfidf_testing_norm)

```

```
199         y_pred = torch.argmax(y_pred, dim=1)
200         acc = accuracy_score(y_test, y_pred)
201         epoch_time = time.time() - epoch_start
202         print(f"{epoch+1:5d} | {avg_loss:.6f} | {acc:.6f} | {epoch_time:.2f}s")
203
204     total_time = time.time() - start_time
205     print(f"\n Training completed in {total_time/60:.2f} minutes")
206     print(f"Final test accuracy: {acc:.4f}")
207
208
209 if __name__ == "__main__":
210     process_data()
```

### 3 Hyperparameter Tuning Results and Analysis

Utilizing the FNN provided by the previous task we attempted to find the learning rate and L2 regularization that produced the best accuracy. The learning rates tested were [0.0001, 0.001, 0.01, 0.1] with L2 regularization's of [0, 0.000001, 0.00001, 0.0001, 0.001, 0.01, 0.1]. Due time and hardware limitations feature sizes of 2500, 5000, 10,000 were run using these learning rates and L2 regularization's for 10 epochs. Once the best learning rate was identified and top two L2 regularization's considered, runs using 20,000 and 50,000 features were made to see if we could gain a better accuracy using larger feature sets. We were able to determine that a learning rate of 0.01 with an L2 regularization of 0.01 performed best for our FNN. Below is a table with accuracies and runtime per epoch for each feature size. We decided to compare accuracies at epoch 8 as it took 396.15 minutes for the 50,000 feature run on our FNN and capped that feature set there.

Table 1: Final Hyperparameter Tuning for FNN: Learning Rate: 0.01, L2-Regularization 0.01

Feature Size	Accuracy (%) at Epoch = 8	Time Range per Epoch (s)
2500	0.8713	11-12 (s)
5000	0.8847	33-38 (s)
10000	0.8899	145-190 (s)
20000	0.8949	560-590 (s)
50000	0.8950	2980-3100 (s)

In Table 1 we can see that as we increase the number of features with the tuned hyperparameter's our accuracy increases. However, the cost per epoch also increase drastically. With enough time and hardware to use the full feature set it may be possible to get over 90% accuracy using our FNN. However, by the time we reached 50,000 features we were at or near hardware capacity for the computers used.

A slight modification of the textbooks logistic regression to use PyTorch and tensors for the dataset was made but the same structure was maintained from the textbook implementation. For the logistic regression it was found that it had the best accuracy at 500 iterations (iterations tested were [10, 100, 500, 1000, 2000]). So using the same learning rate of 0.01 at 500 iterations we have the following table for logistic regression.

Table 2: Logistic Regression: Learning Rate: 0.01, Iterations = 500

Feature Size	Accuracy (%)	Total Time Cost (s)
2500	0.8830	5.4 (s)
5000	0.8891	12.6 (s)
10000	0.8891	25.8 (s)
20000	0.8889	46.2 (s)
50000	0.8816	103.8 (s)

The run time of the logistic regression is several times faster that the FNN. However, the accuracy does not seem to be impacted by the number of features add. It seems to swing in the range of 88 to 89 percent and would infer that the accuracy is capped at 89% for the logistic regression model of this dataset.

## 4 K-Fold Cross Validation Analysis

### 4.0.1 Impact of $k$ -Fold Cross Validation

For the baseline configuration ( $K = 1$ ), the network attains 88.28 % accuracy in 3.20 s. Introducing  $k$ -fold cross validation improves accuracy only marginally—at  $K = 5$ , the model rises to 88.25 % and at  $K = 10$ , it converges near 88.65 %. However, wall-clock time increases almost linearly with  $K$ , as shown in Table 3. This test case demonstrates that although cross validation yields a slightly more reliable estimate of performance, it imposes an immense overhead that is seldom justified once the architecture and hyper-parameters have been fixed.

Table 3: K-Fold Cross-Validation with the Original Model: Time Cost and Mean Accuracy

$K$	Accuracy (%)	Time Cost (s)
1	0.8828	3.20
2	0.8758	5.10
3	0.8790	9.25
4	0.8804	14.90
5	0.8825	19.78
6	0.8815	23.04
7	0.8803	27.34
8	0.8828	32.73
9	0.8837	37.86
10	0.8865	43.92
11	0.8866	49.12
12	0.8867	54.56
13	0.8864	59.20

```
1 def create_model():
2     return nn.Sequential(
3         nn.Linear(INPUT_SIZE, 1024),
4         nn.ReLU(),
5         nn.Linear(1024, 512),
6         nn.ReLU(),
7         nn.Linear(512, 128),
8         nn.ReLU(),
9         nn.Linear(128, 2)
10    )
```

### 4.0.2 Smaller Model vs. Original Model

After experimenting with different builds, The tuned architecture delivers comparable or better accuracies while roughly halving execution time at every  $K$  (e.g., 21.76 s vs. 43.92 s at  $K = 10$ ). This behavior could be due to fewer parameters shortening forward/backward passes, directly cutting runtime. In addition, the lower capacity mitigates over-fitting on the features, allowing validation accuracy to match or even slightly exceed that of the deeper original network.



Table 4: K-Fold Cross-Validation with a Tuned Model: Time Cost and Mean Accuracy

$K$	Accuracy (%)	Time Cost (s)
1	0.8830	2.28
2	0.8758	3.32
3	0.8759	5.50
4	0.8848	7.84
5	0.8840	10.16
6	0.8855	12.49
7	0.8783	14.79
8	0.8857	17.17
9	0.8841	19.42
10	0.8865	21.76
11	0.8868	24.05
12	0.8869	26.14
13	0.8865	31.52

## 5 Training using Dropout Regularization

To add dropout to our FNN, we used PyTorch’s Dropout function to dropout following the last ReLU layer. The dropout probability was set to 0.5 for all models trained. Below is the function that generated the FNN with dropout.

```

1 def gen_single_dropout():
2     print("\n4. _Initializing_neural_network...")
3     net_single_dropout = torch.nn.Sequential(
4         torch.nn.Linear(10000, 5000),
5         torch.nn.ReLU(),
6         torch.nn.Linear(5000, 5000),
7         torch.nn.ReLU(),
8         torch.nn.Linear(5000, 5000),
9         torch.nn.ReLU(),
10        torch.nn.Linear(5000, 2500),
11        torch.nn.ReLU(),
12        torch.nn.Dropout(p=0.5),
13        torch.nn.Linear(2500, 2),
14        torch.nn.Softmax(dim=1))
15
16     return net_single_dropout

```

The same training code used to build the baseline model was used to train the single dropout model. For the set of bagged dropout models, the script below was used to train five models and then ensemble them.

```

1 print("===_Bagging_5_Dropout_Models_Training_===")
2 bag_accs = []
3 bag_times = []
4 for i in range(5):
5     print(f"--_Training_model_{i+1}_--")
6     model = gen_single_dropout()
7     acc, t = train(model, train_loader, X_test, y_test)
8     bag_accs.append(acc)
9     bag_times.append(t)

```

```

10
11
12 bag_acc = [sum(epoch_accs) / len(bag_accs) for epoch_accs in zip(*bag_accs)]
13 total_bag_time = sum(bag_times)
14 print(f"Bagging_total_time:_{total_bag_time/60:.2f}_minutes, _Final_Ensemble_Avg_Acc:_
      {bag_acc[-1]:.4f}\n")

```

Model Version	Training Time (minutes)
Baseline	13.83
Single Dropout	12.91
Bagging (5 Dropout Models)	52.32

Table 5: Comparison of training times for different model configurations to 12 epochs. Max feature input was set to 10000.

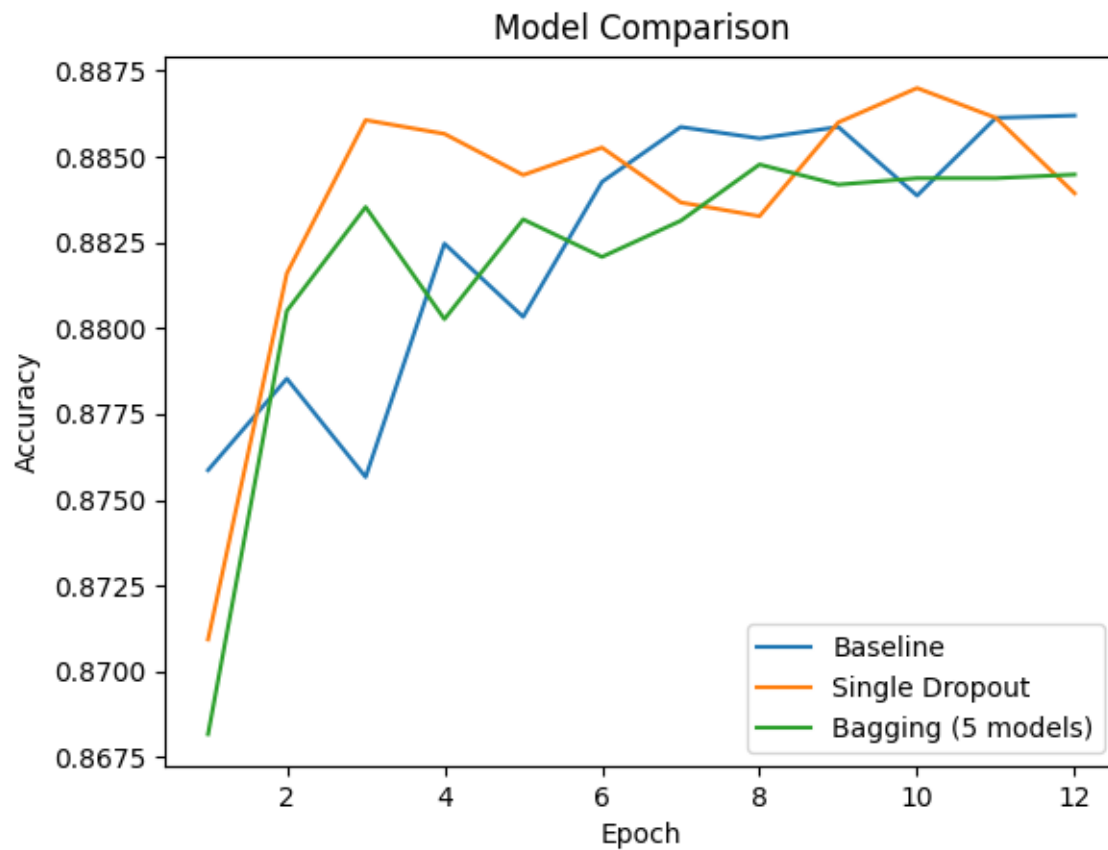


Figure 1: Accuracy over epochs for model variants. Max feature input was set to 10000.

## 6 Team Member Contributions

Below are the contributions of each group member.

**Ellie Larence:** Worked on Task 5, and helped give clarity on Tasks 2, 3, and 4.

**Matthew Lloyd Macias:** Worked on tasks 1, 2, and ran a few supplementary experimental tests for Task 3.

**Jamerson Tenorio:** Worked on Task 3, reformatting the original `generate_data.py` script to be more modular and be able to experiment with variable defined hyperparameters for tuning.

**Leo Wang:** Worked on Task 4, as well as provided some needed fixes to the original `generate_data.py` script.

## References

- [1] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL: <http://www.aclweb.org/anthology/P11-1015>.

[1]