

Technical Report: Final Project DS 5110: Introduction to Data Management and Processing

Team Members: Colin Johnson
Khoury College of Computer Sciences
Data Science Program
`johnson.colin@northeastern.edu`

December 10, 2024

Contents

| | | |
|----------|-------------------------------|----------|
| 1 | Introduction | 2 |
| 2 | Literature Review | 2 |
| 3 | Methodology | 2 |
| 3.1 | Data Collection | 2 |
| 3.2 | Data Preprocessing | 2 |
| 3.3 | Analysis Techniques | 2 |
| 4 | Results | 3 |
| 5 | Discussion | 4 |
| 6 | Conclusion | 4 |
| 7 | References | 5 |
| A | Appendix A: Code | 5 |

1 Introduction

In this project I am investigating the Stanford Movie Dataset which includes a split of positive and negative train and test data. The goal of this project is to make a model that can predict the sentiment score of reviews by using PyTorch Machine Learning.

2 Literature Review

Other projects of the same time use varying types of models. Some use RNN model which is more standard however a lot use CNN too. Models similar to mine found similar accuracy results depending on how well the data is preprocessed.

3 Methodology

I found the RNN models to be lacking in depth when I applied it to my dataset, I would have lower train accuracy. After applying CNN the model fit much better.

3.1 Data Collection

I collected the data by using `os` and extracted the tarfile using the `tarfile` import. I split the data into two separate datasets, one with the training data and one with the testing data, contain both all the negative and positive labels.

3.2 Data Preprocessing

Processed by using `re` to remove HTML labels and punctuations. I lemmatized the data to generalize the words filter specific words that were popular but useless in order to make the runtime faster (words like "movie", "film", etc.). Then I tokenized the reviews to only include important words.

After, I then took every token and matched it with vocab index values to encode words essentially. Words like ['love'] would turn into [200] for example here. I finally padded to the 90th percentile as for the model to work all data has to have equal length.

3.3 Analysis Techniques

Used a CNN model to relate the text to a sentiment score using the label. An example of how this model would work is if you had a tokenized review like (['love'], ['great'], ['happy']). The model then would detect the label is 1 therefore it relates these terms to positive. After, if it sees this more and more then it relates it more to positive.

Parameters

- **vocab_size**: Size of the vocabulary.
- **embedding_dim**: Dimensionality of the word embeddings.
- **conv_config**: Configuration of convolutional layers, including:

- **kernel_sizes**: List of kernel sizes (window sizes for the filters).
- **num_channels**: Number of filters per kernel size.
- **output_size**: Number of output classes (pos, neg)
- **dropout**: Dropout probability to prevent overfitting the padded sequences

4 Results

Results found that the model worked with 85 percent accuracy after 10 epochs. Here

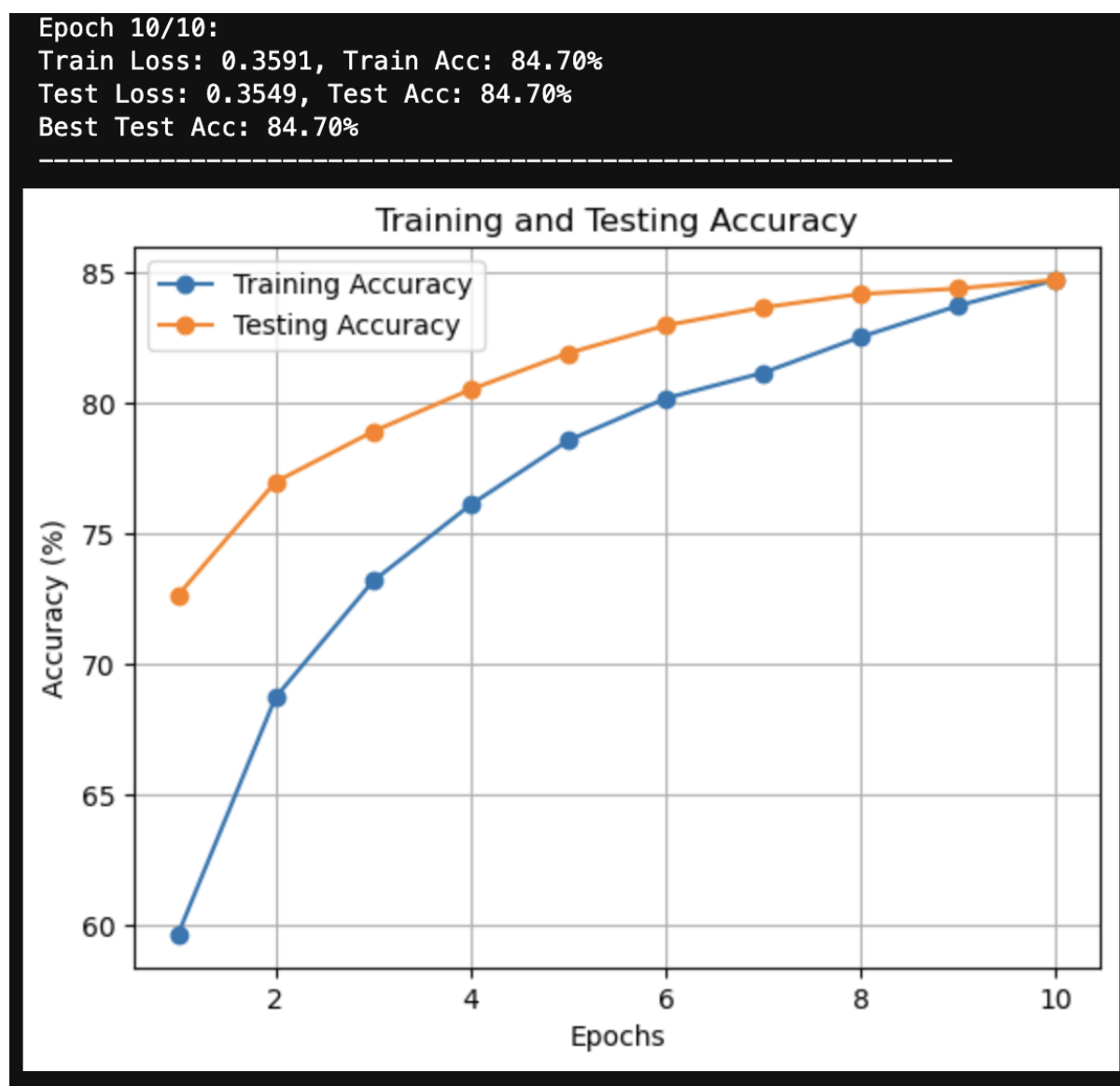


Figure 1: Accuracy after 10 epochs

shows my model is 85 percent accurate. I ran a test case for the model to figure how accurate it was and here they are

```

Review: This movie is a gem that reminds us why we love cinema. From the very first scene, it pulls you into a vivid world of compelling characters, breathtaking visuals, and a narrative that keeps you emotionally invested until the very last frame. The performances are nothing short of extraordinary. The lead actor delivers a career-defining role, capturing every nuance of the character's journey. The supporting cast is equally remarkable, creating a dynamic ensemble that feels authentic and relatable. The director's vision is evident in every scene, blending artistry and storytelling in a way that feels seamless. The cinematography is stunning, with each frame looking like a carefully crafted painting. The music score adds another layer of depth, perfectly complementing the emotional highs and lows of the story. What truly sets this movie apart is its heart. It's a film that resonates on a deeply personal level, leaving you with a sense of hope, wonder, and inspiration. Whether you're a fan of heartfelt dramas, epic romances, or thought-provoking narratives, this film has something for everyone. A must-watch for anyone who loves great storytelling, this movie will stay with you long after the credits roll.
Predicted Sentiment: Positive Review

Review: This movie is a huge disappointment that fails to live up to its promises. From the very first scene, it struggles to capture your attention and quickly becomes tedious. The characters are bland and uninteresting, with performances that lack depth. The plot is predictable and meanders aimlessly, never building any real tension or excitement. The director's vision is absent, and the cinematography feels lazy and uninspired. The music score is generic and adds nothing to the story. Overall, it's a film that fails to connect on any meaningful level, leaving you feeling empty and unsatisfied. A total miss for anyone expecting a compelling cinematic experience.
Predicted Sentiment: Negative Review

Review: I hate this movie so much, this movie is very good at being boring and slow
Predicted Sentiment: Negative Review

Review: This is a beautifully written movie!
Predicted Sentiment: Positive Review

Review: This movie is an absolute masterpiece, delivering a cinematic experience that is both deeply moving and visually stunning. From the very first scene, it captivates the audience with its rich storytelling and compelling characters. The performances are nothing short of spectacular, with the cast delivering their roles with remarkable depth and authenticity. The lead actor gives a tour-de-force performance, showcasing an impressive range of emotions that make their character unforgettable. The director's vision is evident in every frame, crafting a story that is as thought-provoking as it is emotionally resonant. The cinematography is breathtaking, with beautifully composed shots that feel like works of art. The musical score is equally powerful, perfectly complementing the highs and lows of the narrative and drawing the audience even deeper into the story. What sets this film apart is its ability to balance grand, epic storytelling with intimate, personal moments that strike a universal chord. It's a rare gem that combines artistry and entertainment, leaving you with a sense of awe and inspiration. This is a movie that will stay with you long after the credits roll, reminding you of the magic of great cinema. Truly a must-watch for film lovers everywhere!
Predicted Sentiment: Positive Review

Review: This movie is a complete letdown, failing to deliver on even the most basic expectations of a good film. The story is uninspired and painfully predictable, offering no surprises or moments of genuine intrigue. The characters are one-dimensional, making it impossible to care about their journeys or outcomes. The performances feel flat and lack any emotional depth, as though the cast themselves weren't invested in the material. The direction is sloppy, with disjointed pacing that makes the film drag unbearably. Visually, the cinematography is bland and uninspired, lacking creativity or style. The musical score does little to enhance the experience, feeling generic and forgettable. Overall, this movie is an exercise in mediocrity, leaving no lasting impression other than frustration. It's a soulless production that feels more like a chore than entertainment. A complete waste of time for anyone seeking a compelling or enjoyable cinematic experience.
Predicted Sentiment: Negative Review

Review 1: Positive Review (Confidence: 0.96)
Review 2: Negative Review (Confidence: 0.98)
Review 3: Negative Review (Confidence: 0.92)
Review 4: Positive Review (Confidence: 0.67)
Review 5: Positive Review (Confidence: 0.99)
Review 6: Negative Review (Confidence: 0.98)

```

Figure 2: Test case for model

5 Discussion

The results show there is some work in progress still with my model. For example, notice how review 4 (the shortest one) has confidence value of .67. This implies shorter reviews are less accurate which is why the accuracy for the model is 85 percent. The other figures represent the model working well especially for longer more informed reviews. Compared to other works well others have better accuracy due to stronger preprocessing so if I had the time to change any part it would be the preprocessing and make it clearer which review is positive or not for shorter reviews. However, there is also just error with the model in finding reviews here. This is because there is sometimes so little info it will be very hard to fine-tune to the point that you can get it perfect every time.

6 Conclusion

Found that CNN works well for Sentiment Analysis compared to RNN models at least for my dataset. Model accuracy is highly dataset related as despite there being more work for tokenization, if the dataset is not large enough you can get poor performance due to lack there of references for the model to work with. This is also a limitation as smaller datasets will not be able to get all the nuances of reviews that come their way, causing for lower accuracy. For future research, I can implement a neutral section where

I investigate with neutral reviews too, hopefully making for better accuracy due to more categories for the model to reference when building. Dealing with smaller less nuanced reviews as well, with more research and time it is possible to make the model be able to investigate vague reviews as well.

7 References

References

here

here

A Appendix A: Code

```

1 def tokenize_with_vocab(text, vocab):
2     '''
3     Tokenizing the reviews, removing common and useless words, etc.
4     '''
5     words_remove = ['movie', 'film', 'director', 'plays', 'horror', '
6     comedy', 'watching', 'seen', 'people', 'guy', 's', 'time', 'second',
7     'book']
8     vocab = filter_vocab(vocab, words_remove)
9     tokenized_reviews = []
10    for doc in nlp.pipe(text, batch_size=1000, disable=["parser", "ner"
11    , "tagger"]):
12        cleaned_text = remove_html_tags(doc.text)
13        cleaned_text = remove_punctuation(doc.text)
14        lemmas = spacy_lemmatization(cleaned_text)
15        tokens = [
16            lemma.lower() for lemma in lemmas
17            if lemma.lower() in vocab and lemma.lower() not in nlp.
18            Defaults.stop_words
19        ]
20        tokenized_reviews.append(tokens)
21
22    return tokenized_reviews

```

Listing 1: Preprocessing

```

1 class TextCNN(nn.Module):
2     def __init__(self, vocab_size, embedding_dim, conv_config,
3     output_size, dropout=0.5):
4         super(TextCNN, self).__init__()
5
6         self.conv_config = conv_config
7         self.output_size = output_size
8         self.dropout_p = dropout
9
10        self.embedding = nn.Embedding(vocab_size, embedding_dim)
11        #configuring the hidden layers, extracting features from the
12        input nodes (reviews)
13        self.convolutions = nn.ModuleList([
14            nn.Sequential(

```

```

14         nn.Conv1d(embedding_dim, self.conv_config['num_channels
15         ], kernel_size=kernel),
16         nn.ReLU(),
17         nn.AdaptiveMaxPool1d(1)
18     )
19     for kernel in self.conv_config['kernel_sizes']
20 ]])
21
22     self.dropout = nn.Dropout(self.dropout_p)
23     self.linear = nn.Linear(
24         self.conv_config['num_channels'] * len(self.conv_config['
25         kernel_sizes']),
26         self.output_size
27     )
28
29     def forward(self, input_seq):
30         '''
31         Forward pass for the TextCNN model.
32         Returns the log probability with shape of output and batch size
33         '''
34
35         emb_out = self.embedding(input_seq).permute(0, 2, 1)
36
37         conv_out = [conv(emb_out).squeeze(2) for conv in self.
38         convolutions]
39
40         concat_out = torch.cat(conv_out, dim=1)
41
42         concat_out = self.dropout(concat_out)
43         out = self.linear(concat_out)
44
45         return F.log_softmax(out, dim=-1)
46
47 class SentimentDataset(Dataset):
48     def __init__(self, sequences, labels):
49         self.sequences = sequences
50         self.labels = torch.tensor(labels, dtype=torch.long)
51
52     def __len__(self):
53         return len(self.sequences)
54
55     def __getitem__(self, idx):
56         return self.sequences[idx], self.labels[idx]

```

Listing 2: CNN Model

```

1 def train_model(model, train_loader, test_loader, criterion, optimizer,
2   num_epochs=10):
3     device='cpu'
4     model = model.to(device)
5     train_losses = []
6     test_losses = []
7     train_accuracies = []
8     test_accuracies = []
9     best_acc = 0.0
10
11     for epoch in range(num_epochs):
12         model.train()

```

```

12     running_loss = 0.0
13     correct = 0
14     total = 0
15
16     for inputs, labels in train_loader:
17         inputs, labels = inputs.to(device), labels.to(device)
18         optimizer.zero_grad()
19         outputs = model(inputs)
20         #comparing how much you lose compared to output target
values
21         loss = criterion(outputs, labels)
22         loss.backward()
23         optimizer.step()
24
25         running_loss += loss.item()
26         _, predicted = outputs.max(1)
27         total += labels.size(0)
28         correct += predicted.eq(labels).sum().item()
29
30
31     epoch_loss = running_loss / len(train_loader)
32     epoch_acc = 100. * correct / total
33     train_losses.append(epoch_loss)
34     train_accuracies.append(epoch_acc)
35
36
37     model.eval()
38     test_loss = 0.0
39     correct = 0
40     total = 0
41
42     with torch.no_grad():
43         for inputs, labels in test_loader:
44             inputs, labels = inputs.to(device), labels.to(device)
45             outputs = model(inputs)
46             loss = criterion(outputs, labels)
47             #collect test loss & test values
48             test_loss += loss.item()
49             _, predicted = outputs.max(1)
50             total += labels.size(0)
51             correct += predicted.eq(labels).sum().item()
52
53
54     test_loss = test_loss / len(test_loader)
55     test_acc = 100. * correct / total
56     test_losses.append(test_loss)
57     test_accuracies.append(test_acc)
58
59     if test_acc > best_acc:
60         best_acc = test_acc
61         torch.save(model.state_dict(), 'best_model.pth')
62
63     print(f'Epoch {epoch+1}/{num_epochs}:')
64     print(f'Train Loss: {epoch_loss:.4f}, Train Acc: {epoch_acc:.2f
}%')
65     print(f'Test Loss: {test_loss:.4f}, Test Acc: {test_acc:.2f}%')
66     print(f'Best Test Acc: {best_acc:.2f}%')
67     print('-' * 60)

```

```
68     epochs = range(1, len(train_accuracies) + 1)
69     '''
70     Plot using matplotlib to display the change in accuracy over
multiple epoch's
71     '''
72     plt.plot(epochs, train_accuracies, label='Training Accuracy',
marker='o')
73     plt.plot(epochs, test_accuracies, label='Testing Accuracy',
marker='o')
74     plt.xlabel('Epochs')
75     plt.ylabel('Accuracy (%)')
76     plt.title('Training and Testing Accuracy')
77     plt.legend()
78     plt.grid(True)
79     plt.show()
```

Listing 3: Training Model

```
1 review = ["This movie is a gem that reminds us why we love cinema. From
the very first scene, it pulls you into a vivid world of compelling
characters, breathtaking visuals, and a narrative that keeps you
emotionally invested until the very last frame. \
2 The performances are nothing short of extraordinary. The lead actor
delivers a career-defining role, capturing every nuance of the
character's journey. The supporting cast is equally remarkable,
creating a dynamic ensemble that feels authentic and relatable. \
3 The director's vision is evident in every scene, blending artistry and
storytelling in a way that feels seamless. The cinematography is
stunning, with each frame looking like a carefully crafted painting.
The music score adds another layer of depth, perfectly \
4 complementing the emotional highs and lows of the story. What truly
sets this movie apart is its heart. It's a film that resonates on a
deeply personal level, leaving you with a sense of hope, wonder, and
inspiration. Whether you're a fan of heartfelt dramas, \
5 epic romances, or thought-provoking narratives, this film has something
for everyone. A must-watch for anyone who loves great storytelling,
this movie will stay with you long after the credits roll.",
6 "This movie is a huge disappointment that fails to live up to \
7 its promises. From the very first scene, it struggles to capture your
attention and quickly becomes tedious. The characters are bland and
uninteresting, with performances that lack depth. The plot is
predictable and meanders aimlessly, never building any real tension
or \
8 excitement. The director's vision is absent, and the cinematography
feels lazy and uninspired. The music score is generic and adds
nothing to the story. Overall, it's a film that fails to connect on
any meaningful level, leaving you feeling empty and unsatisfied. \
9 A total miss for anyone expecting a compelling cinematic experience.",
10 "I hate this movie so much, this movie is very good at being boring and
slow",
11 "This is a beautifully written movie!",
12 "This movie is an absolute masterpiece, delivering a cinematic
experience that is both deeply moving and visually stunning. From
the very first scene, it captivates the audience with its rich
storytelling and compelling characters. The performances are nothing
short of spectacular, with the cast delivering their roles with
remarkable depth and authenticity. The lead actor gives a tour-de-
force performance, showcasing an impressive range of emotions that
make their character unforgettable. \
```



```

13 The director's vision is evident in every frame, crafting a story that
    is as thought-provoking as it is emotionally resonant. The
    cinematography is breathtaking, with beautifully composed shots that
    feel like works of art. The musical score is equally powerful,
    perfectly complementing the highs and lows of the narrative and
    drawing the audience even deeper into the story. \
14 What sets this film apart is its ability to balance grand, epic
    storytelling with intimate, personal moments that strike a universal
    chord. It's a rare gem that combines artistry and entertainment,
    leaving you with a sense of awe and inspiration. This is a movie
    that will stay with you long after the credits roll, reminding you
    of the magic of great cinema. Truly a must-watch for film lovers
    everywhere!",
15 "This movie is a complete letdown, failing to deliver on even the most
    basic expectations of a good film. The story is uninspired and
    painfully predictable, offering no surprises or moments of genuine
    intrigue. The characters are one-dimensional, making it impossible
    to care about their journeys or outcomes. \
16 The performances feel flat and lack any emotional depth, as though the
    cast themselves weren't invested in the material. The direction is
    sloppy, with disjointed pacing that makes the film drag unbearably.
    Visually, the cinematography is bland and uninspired, lacking
    creativity or style. The musical score does little to enhance the
    experience, feeling generic and forgettable. \
17 Overall, this movie is an exercise in mediocrity, leaving no lasting
    impression other than frustration. It's a soulless production that
    feels more like a chore than entertainment. A complete waste of time
    for anyone seeking a compelling or enjoyable cinematic experience."
    ]
18 '''
19 Test case
20 '''
21 df = pd.DataFrame(review, columns=['review'])
22
23 df['tokenized_review'] = tokenize_with_vocab(df['review'], vocab)
24 print(df['tokenized_review'])
25 df['review_ids'] = df['tokenized_review'].apply(lambda tokens:
    tokenize_id(tokens, vocab))
26 df['review_length'] = df['tokenized_review'].apply(len)
27 sequence_df = [torch.tensor(ids, dtype=torch.long) for ids in df['
    review_ids']]
28 max_length = int(df['review_length'].max())
29 padded = pad_sequences(sequence_df, max_length)
30
31 model.eval()
32 with torch.no_grad():
33     output = model(padded)
34
35     predicted_class = torch.argmax(output, dim=1).tolist()
36
37 label_map = {0: "Negative Review", 1: "Positive Review"}
38 predicted_labels = [label_map[p] for p in predicted_class]
39
40 for i, review in enumerate(review):
41     print(f"Review: {review}")
42     print(f"Predicted Sentiment: {predicted_labels[i]}\n")
43 probabilities = torch.softmax(output, dim=1)
44 for i, prob in enumerate(probabilities):

```

```
45     if prob[0] > prob[1]:
46         sentiment = "Negative Review"
47     else:
48         sentiment = "Positive Review"
49     print(f"Review {i + 1}: {sentiment} (Confidence: {prob.max().item():.2f})")
```

Listing 4: Test Cases