

```
In [1]: import numpy as np
import pandas as pd
import torch
import torch.nn.functional as F
import torchtext
import time
import random
import pandas as pd
import spacy
```

Import Required Libraries & Data Loading

```
In [2]: #importing the training data
df=pd.read_csv('IMDB Dataset.csv')
print(df.shape)
df.head(10)
```

```
Out[2]:
```

	review	sentiment
0	One of the other reviewers has mentioned that...	positive
1	A wonderful little production. The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Matte's 'Love in the Time of Money' is...	positive
5	Probably my all-time favorite movie, a story o...	positive
6	I sure would like to see a resurrection of a u...	positive
7	This show was an amazing, fresh & innovativ...	negative
8	Encouraged by the positive comments about th...	negative
9	If you like original gut wrenching laughter yo...	positive

Data Preparation

```
In [3]: """
sentiment : 0 = negative, 1 = positive
use the following to get the sentiment of a sentence :
sentiment = 0 if sentiment is negative else 1

use np.where to get the sentiment of a sentence :
"""
df['sentiment'] = np.where(df['sentiment'] == 'positive', 1, 0)

In [4]: df.head()
```

```
Out[4]:
```

	review	sentiment
0	One of the other reviewers has mentioned that ...	1
1	A wonderful little production. The...	1
2	I thought this was a wonderful way to spend ti...	1
3	Basically there's a family where a little boy ...	0
4	Petter Matte's 'Love in the Time of Money' is...	1

```
In [5]: df.columns = ['TEXT_COLUMN_NAME', 'LABEL_COLUMN_NAME']

In [6]: """
Load the spacy model and load the English language model from https://spacy.io/usage/models
"""
nlp = spacy.load('en_core_web_sm')
```

```
In [7]: # general Settings

RANDOM_SEED = 123
torch.manual_seed(RANDOM_SEED)

VOCABULARY_SIZE = 20000
LEARNING_RATE = 0.002 ### ADD YOUR LEARNING RATE HERE ###
BATCH_SIZE = 50 ### ADD YOUR BATCH SIZE HERE ###
NUM_EPOCHS = 25 ### ADD YOUR NUMBER OF EPOCHS HERE ###
DEVICE = torch.device('cuda' if torch.cuda.is_available() else 'cpu')

EMBEDDING_DIM = 100 ### ADD YOUR EMBEDDING DIMENSION HERE ###
HIDDEN_DIM = 100 ### ADD YOUR HIDDEN DIMENSION HERE ###
NUM_CLASSES = 2
```

Text & label Preparation

```
In [8]: # Define feature processing
"""
Define the fields for the data.
"""
TEXT = torchtext.legacy.data.Field(tokenize = 'spacy', tokenizer_language = 'en_core_web_sm')
```

```
In [9]: # Define Label processing
LABEL = torchtext.legacy.data.LabelField(dtype = torch.long)
```

```
In [10]: """
Define the fields for the data.
"""

df.to_csv('moviedata.csv', index = None)
df = pd.read_csv('moviedata.csv')
df.head()
```

```
Out[10]:
```

	TEXT_COLUMN_NAME	LABEL_COLUMN_NAME
0	One of the other reviewers has mentioned that ...	1
1	A wonderful little production. The...	1
2	I thought this was a wonderful way to spend ti...	1
3	Basically there's a family where a little boy ...	0
4	Petter Matte's 'Love in the Time of Money' is...	1

```
In [11]: # process the dataset

fields = [('TEXT_COLUMN_NAME', TEXT), ('LABEL_COLUMN_NAME', LABEL)]

dataset = torchtext.legacy.data.TabularDataset({
    path = "moviedata.csv", ### ADD YOUR DATASET PATH HERE ###
    format = "CSV", ### ADD YOUR DATASET FORMAT HERE ###
    skip_header = True, ### ADD YOUR SKIP HEADER HERE ###
    fields = fields ### ADD YOUR FIELDS HERE ###
})
```

Data Split

```
In [12]: # Split dataset into train and test set

train_data, test_data = dataset.split(split_ratio = [0.8, 0.2], random_state =
random.seed(RANDOM_SEED))

print('Length of train data', len(train_data))
print('Length of test data', len(test_data))

Length of train data 40000
Length of test data 10000

In [13]: train_data, val_data = train_data.split(split_ratio = [0.85, 0.15], random_state =
random.seed(RANDOM_SEED))

print('Length of train data', len(train_data))
print('Length of valid data', len(val_data))

Length of train data 34000
Length of valid data 6000
```

Data Observation after Tokenization

```
In [14]: # Look at first training example

print(vars(train_data.examples[2009]))

{'TEXT_COLUMN_NAME': ['There', 'was', 'a', 'lot', 'of', 'hype', 'of', 'this', 'movie', 'and', 'the', 'commercials', 'made', 'it', 'seem', 'like', 'it', 'would', 'be', 'great', 'I', 'saw', 'it', 'like', 'Bling', 'It', 'o n', '2', 'I', 'Bling', 'It', 'On', '3', 'shamed', 'glory', 'of', 'the', 'original', 'Bling', 'It', 'on', 'I', 'There', 'is', 'shameless', 'stereotyping', 'throughout', 'the', 'film', 'I', 'The', 'lines', 'given', 'to', 't he', 'actors', 'were', 'humiliating', 'for', 'all', 'the', 'races', 'involved', 'In', 'the', 'film', 'I', 'The e', 'performance', 'of', 'Hayden', 'Panettiere', 'was', 'sub', 'I', 'par', 'both', 'in', 'terms', 'of', 'actin g', 'and', 'cheerleading', 'I', 'There', 'were', 'several', 'scenes', 'in', 'which', 'I', 'literally', 'cringe d', 'because', 'I', 'was', 'embarrassed', 'for', 'the', 'cast', 'because', 'the', 'scene', 'I', 'lines', 'I', 'plot', 'I', 'etc', 'I', 'were', 'just', 'so', 'stupid', 'I', 'My', 'recommendation', 'to', 'the', 'makers', 'to f', 'any', 'future', 'Bling', 'It', 'On', 'film', 'is', 'that', 'you', 'should', 'hire', 'good', 'cheerleader s', 'and', 'teach', 'them', 'to', 'act', 'because', 'the', 'I', 'acting', 'I', 'of', 'the', 'cast', 'was', 'horrendous', 'and', 'thei', 'lack', 'of', 'cheerleading', 'ability', 'made', 'them', 'completely', 'useless', 't o', 'the', 'film', 'I', 'Only', 'great', 'character', 'I', 'Rirresha', 'I', 'LABEL_COLUMN_NAME': 0}]

In [15]: # Build Vocabulary

TEXT.build_vocab(train_data, max_size = VOCABULARY_SIZE)
LABEL.build_vocab(train_data)

print(f'vocabulary size: {len(TEXT.vocab)}')
print(f'Label Size: {len(LABEL.vocab)}')
```

vocabulary size: 20002
Label Size: 2

2 extra value in vocabulary is because added (unknown) and (padding)

```
In [16]: # Print the most common words: Use the most common method of the TEXT vocabulary
most_common_words = TEXT.vocab.freqs.most_common()
print(most_common_words[:20])

[('the', 390972), ('I', 369444), ('.', 318509), ('a', 210502), ('and', 210008), ('of', 194659), ('to', 180163), ('is', 145895), ('in', 118266), ('I', 105681), ('it', 103588), ('that', 93995), ('"', 85535), ('"s', 83149), ('this', 81775), ('-', 71249), ('/s>cb', 68787), ('was', 67372), ('as', 57734), ('movie', 57572)]
```

```
In [17]: # Token corresponding to first 10 Indices

print(TEXT.vocab.itos[:20]) #itos = Integer to string

['<unk>', '<pad>', 'the', 'I', 'a', 'and', 'of', 'to', 'is', 'in', 'it', 'it', 'that', '', '"s', 'this', 'I', 's>cb', 'was']
```

Data Preparation for Batch wise Implementation

```
In [18]: # Define DataLoader

train_loader, valid_loader, test_loader = torchtext.legacy.data.BucketIterator.splits(
(train_data, val_data, test_data), ### ADD YOUR SPLIT DATA HERE (Make sure you add it in a tuple) ###
    batch_size = BATCH_SIZE, ### ADD YOUR BATCH SIZE HERE ###
    sort_within_batch = True, ### ADD YOUR SORT WITHIN BATCH HERE ###
    sort_key = lambda x : len(x.TEXT_COLUMN_NAME),
    device = DEVICE
)
```

```
In [19]: # Testing the divisors (note that the number of rows depends on the longest document in the
respective batch):

print('Train')
for batch in train_loader:
    print(f'Text matrix size: {batch.TEXT_COLUMN_NAME.size()}')
    print(f'Target vector size: {batch.LABEL_COLUMN_NAME.size()}')
    break

print('\nValid:')
for batch in valid_loader:
    print(f'Text matrix size: {batch.TEXT_COLUMN_NAME.size()}')
    print(f'Target vector size: {batch.LABEL_COLUMN_NAME.size()}')
    break

print('\nTest:')
for batch in test_loader:
    print(f'Text matrix size: {batch.TEXT_COLUMN_NAME.size()}')
    print(f'Target vector size: {batch.LABEL_COLUMN_NAME.size()}')
    break

Train
Text matrix size: torch.Size([250, 50])
Target vector size: torch.Size([50])

Valid:
Text matrix size: torch.Size([47, 50])
Target vector size: torch.Size([50])

Test:
Text matrix size: torch.Size([41, 50])
Target vector size: torch.Size([50])
```

Model Building

```
In [49]: import torch.nn as nn

class RNN(nn.Module):

    def __init__(self, input_dim, embedding_dim, hidden_dim, output_dim):
        super().__init__()
        ### ADD YOUR CODE HERE ###
        self.hidden_dim = hidden_dim

        self.embed = nn.Embedding(input_dim, embedding_dim)
        self.rnn = nn.LSTM(embedding_dim, hidden_dim, num_layers=1)
        self.fc = nn.Linear(hidden_dim, output_dim)
        ### END YOUR CODE ###

    def forward(self, text):
        ### ADD YOUR CODE HERE ###
        # text dim: [sentence length, batch size]
        # embedded dim: [sentence length, batch size, embedding dim]
        embedded = self.embed(text)
        output, (hidden, cell) = self.rnn(embedded)
        # hidden dim: [1, batch size, hidden dim]
        hidden.squeeze(0)
        # hidden dim: [batch size, hidden dim]
        output = self.fc(hidden)
        ### END YOUR CODE ###
        #reshape output to [batch size, output_dim]
        return output.permute(1, 2, 0).squeeze(2)
```

```
In [50]: torch.manual_seed(RANDOM_SEED)

model = RNN(input_dim=len(TEXT.vocab), ### ADD YOUR INPUT DIM HERE. This can be the length of
your vocabulary or the embedding dim ###
            embedding_dim= EMBEDDING_DIM, ### ADD YOUR EMBEDDING DIM HERE ###
            hidden_dim= HIDDEN_DIM, ### ADD YOUR HIDDEN DIM HERE ###
            output_dim= NUM_CLASSES ### ADD NUMBER OF CLASSES HERE ###
)

model = model.to(DEVICE)

loss_fn = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=LEARNING_RATE) ### ADD YOUR OPTIMIZER HERE
###
```

Define Accuracy

```
In [51]: def compute_accuracy(model, data_loader, device):

    with torch.no_grad():

        correct_pred, num_examples = 0, 0

        for i, (features, targets) in enumerate(data_loader):

            features = features.to(device)
            targets = targets.float().to(device)

            logits = model(features)
            _, predicted_labels = torch.max(logits, 1)

            num_examples += targets.size(0)
            correct_pred += (predicted_labels == targets).sum()

        return correct_pred.float()/num_examples * 100
```

Model Run

```
In [52]: start_time = time.time()

for epoch in range(NUM_EPOCHS):
    model.train()

    for batch_idx, batch_data in enumerate(train_loader):
        text = batch_data.TEXT_COLUMN_NAME.to(DEVICE)
        labels = batch_data.LABEL_COLUMN_NAME.to(DEVICE)

        ### FORWARD AND BACK PROP
        model.zero_grad()
        output = model(text)
        loss = loss_fn(output, labels)
        loss.backward()

        ### UPDATE MODEL PARAMETERS
        optimizer.step()

        ### LOGGING
        if not batch_idx % 50:
            print(f'Epoch: {epoch+1:03d}/{NUM_EPOCHS:03d} | '
                  f'Batch {batch_idx:03d}/{len(train_loader):03d} | '
                  f'Loss: {loss:.4f}')

    with torch.set_grad_enabled(False):
        print(f'training accuracy: '
              f'(compute_accuracy(model, train_loader, DEVICE):.2f)%'
              f'invalid accuracy: '
              f'(compute_accuracy(model, valid_loader, DEVICE):.2f)%')

    print(f'Time elapsed: {(time.time() - start_time)/60:.2f} min')

print(f'Total Training Time: {(time.time() - start_time)/60:.2f} min')
print(f'Test accuracy: (compute_accuracy(model, test_loader, DEVICE):.2f)%')
```


[illegible]

```

    return predic

print('Probabilit
predict sentiment

```