

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import random as rand
5 from string import punctuation
6 from collections import Counter
7 import torch
8 from torch.utils.data import DataLoader, TensorDataset
9 import torch.nn as nn

1 # read data
2 data = pd.read_csv("https://raw.githubusercontent.com/AyeEyeTwoFive/IMDB/master/IMDB%20Dat
3 reviews = data.iloc[:,0]
4 labels = data.iloc[:,2]
5 reviews = reviews.tolist()
6 labels = labels.tolist()
7 reviews = '\n'.join(reviews)
8 # Convert labels to ints
9 labels = [1 if label == 'positive' else 0 for label in labels]
10 labels = np.array(labels)

1 # Process Data
2 reviews = reviews.lower() # convert to lower case
3 allrev = ''.join([c for c in reviews if c not in punctuation]) # remove punctuation
4 revsplit = allrev.split('\n')
5 print ('Number of reviews :', len(revsplit))
6 allrev2 = ' '.join(revsplit)
7 words = allrev2.split() # list words
8 counts = Counter(words) # Count word frequencies
9 total_words = len(words)
10 sorted_words = counts.most_common(total_words)
11 print (counts)
12 # mapping words to int ordered by frequency
13 vocab_to_int = {w:i+1 for i, (w,c) in enumerate(sorted_words)}
14 # Encode words as ints
15 revints = []
16 for review in revsplit:
17     r = [vocab_to_int[w] for w in review.split()]
18     revints.append(r)
19 print (revints[0:3])
20 %matplotlib inline
21 revslen = [len(x) for x in revints]
22 pd.Series(revslen).hist()
23 plt.title('Distribution of Review Lengths')
24 plt.show()
25 pd.Series(revslen).describe()
26
27 #Remove outlier short or long reviews
28 revints = [revints[i] for i, l in enumerate(revslen) if l>0 ]
```

```
5 batch_size = 50
6 train_loader = DataLoader(train_data, shuffle=True, batch_size=batch_size)
7 valid_loader = DataLoader(valid_data, shuffle=True, batch_size=batch_size)
8 test_loader = DataLoader(test_data, shuffle=True, batch_size=batch_size)

1 ### Define LSTM Model
2
3 class ReviewLSTM(nn.Module):
4
5     def __init__(self, vocab_size, output_size, embedding_dim, hidden_dim, n_layers, drop_
6         """
7         Initialize the model by setting up the layers.
8         """
9         super().__init__()
10
11         self.output_size = output_size
12         self.n_layers = n_layers
13         self.hidden_dim = hidden_dim
14
15         # embedding and LSTM
16         self.embedding = nn.Embedding(vocab_size, embedding_dim)
17         self.lstm = nn.LSTM(embedding_dim, hidden_dim, n_layers,
18                             dropout=drop_prob, batch_first=True)
19
20         # dropout
21         self.dropout = nn.Dropout(0.3)
22
23         # linear and sigmoid
24         self.fc = nn.Linear(hidden_dim, output_size)
25         self.sig = nn.Sigmoid()
26
27
28     def forward(self, x, hidden):
29         """
30         Perform a forward pass of our model on some input and hidden state.
31         """
32         batch_size = x.size(0)
33
34
35         embeds = self.embedding(x)
36         lstm_out, hidden = self.lstm(embeds, hidden)
37
38
39         lstm_out = lstm_out.contiguous().view(-1, self.hidden_dim)
40
41         # fully-connected layer
42         out = self.dropout(lstm_out)
43         out = self.fc(out)
44         # sigmoid
45         sig_out = self.sig(out)
46
```

```
18
19 h = net.init_hidden(batch_size)
20
21
22 for inputs, labels in train_loader:
23     counter += 1
24
25     h = tuple([each.data for each in h])
26
27
28     net.zero_grad()
29
30
31     inputs = inputs.type(torch.LongTensor)
32     output, h = net(inputs, h)
33
34     loss = criterion(output.squeeze(), labels.float())
35     loss.backward()
36
37     nn.utils.clip_grad_norm_(net.parameters(), clip)
38     optimizer.step()
39
40     # calc loss
41     if counter % print_every == 0:
42
43         val_h = net.init_hidden(batch_size)
44         val_losses = []
45         net.eval()
46         for inputs, labels in valid_loader:
47
48
49             val_h = tuple([each.data for each in val_h])
50
51             inputs = inputs.type(torch.LongTensor)
52             output, val_h = net(inputs, val_h)
53             val_loss = criterion(output.squeeze(), labels.float())
54
55             val_losses.append(val_loss.item())
56
57         net.train()
58         print("Epoch: {}/{}...".format(e+1, epochs),
59               "Step: {}...".format(counter),
60               "Loss: {:.6f}...".format(loss.item()),
61               "Val Loss: {:.6f}".format(np.mean(val_losses)))
```




```

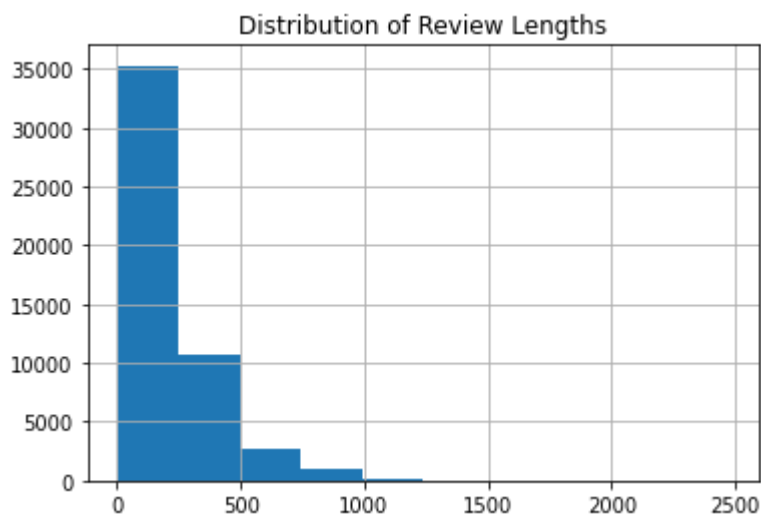
29 labels = [labels[i] for i, l in enumerate(revslen) if l > 0 ]
30
31 # Pad or truncate reviews to achieve consistent length
32 seq_length=200
33 features = np.zeros((len(revints), seq_length), dtype = int)
34
35 for i, review in enumerate(revints):
36     review_len = len(review)
37     if review_len <= seq_length:
38         zeroes = list(np.zeros(seq_length-review_len))
39         new = zeroes+review
40     elif review_len > seq_length:
41         new = review[0:seq_length]
42     features[i,:] = np.array(new)

```



Number of reviews : 50000

Counter({'the': 663815, 'and': 320663, 'a': 320517, 'of': 288382, 'to': 266773, 'is': 21
[1030, 479, 4558, 732, 2375, 1030, 14744, 1376, 2, 3893, 8256, 1505, 1817, 16739, 2, 13



```

1 # 80% train, 10% validation, 10% test splits
2 split_frac = 0.8
3 len_feat = len(features)
4 train_x = features[0:int(split_frac*len_feat)]
5 train_y = labels[0:int(split_frac*len_feat)]
6 remaining_x = features[int(split_frac*len_feat):]
7 remaining_y = labels[int(split_frac*len_feat):]
8 valid_x = remaining_x[0:int(len(remaining_x)*0.5)]
9 valid_y = remaining_y[0:int(len(remaining_y)*0.5)]
10 test_x = remaining_x[int(len(remaining_x)*0.5):]
11 test_y = remaining_y[int(len(remaining_y)*0.5):]

```

1 # Data Loading

```

2 train_data = TensorDataset(torch.FloatTensor(train_x), torch.FloatTensor(train_y))
3 valid_data = TensorDataset(torch.FloatTensor(valid_x), torch.FloatTensor(valid_y))
4 test_data = TensorDataset(torch.FloatTensor(test_x), torch.FloatTensor(test_y))

```

```

47         # reshape
48         sig_out = sig_out.view(batch_size, -1)
49         sig_out = sig_out[:, -1] # get last batch of labels
50
51
52         return sig_out, hidden
53
54
55     def init_hidden(self, batch_size):
56         weight = next(self.parameters()).data
57
58
59         hidden = (weight.new(self.n_layers, batch_size, self.hidden_dim).zero_(),
60                  weight.new(self.n_layers, batch_size, self.hidden_dim).zero_())
61
62         return hidden

```

```

1 # Instantiate and train model
2 vocab_size = len(vocab_to_int)+1 # +1 for the 0 padding
3 output_size = 1
4 embedding_dim = 400
5 hidden_dim = 256
6 n_layers = 2
7 net = ReviewLSTM(vocab_size, output_size, embedding_dim, hidden_dim, n_layers)
8 print(net)

```

```

📄 ReviewLSTM(
  (embedding): Embedding(181686, 400)
  (lstm): LSTM(400, 256, num_layers=2, batch_first=True, dropout=0.5)
  (dropout): Dropout(p=0.3, inplace=False)
  (fc): Linear(in_features=256, out_features=1, bias=True)
  (sig): Sigmoid()
)

```

```

1 # define loss/opt
2 lr=0.001
3
4 criterion = nn.BCELoss()
5 optimizer = torch.optim.Adam(net.parameters(), lr=lr)
6
7
8 # train params
9 epochs = 1
10 counter = 0
11 print_every = 100
12 clip=5
13
14
15 net.train()
16
17 for e in range(epochs):

```

Epoch: 1/1... Step: 100... Loss: 0.607990... Val Loss: 0.584549
 Epoch: 1/1... Step: 200... Loss: 0.600504... Val Loss: 0.600402

```

1 # Get test data stats
2
3 test_losses = []
4 num_correct = 0
5
6 h = net.init_hidden(batch_size)
7
8 net.eval()
9
10 for inputs, labels in test_loader:
11
12     h = tuple([each.data for each in h])
13
14
15     # predict
16     inputs = inputs.type(torch.LongTensor)
17     output, h = net(inputs, h)
18
19     # calc loss
20     test_loss = criterion(output.squeeze(), labels.float())
21     test_losses.append(test_loss.item())
22
23     # threshold probabilities
24     pred = torch.round(output.squeeze())
25
26     # test predictions
27     correct_tensor = pred.eq(labels.float().view_as(pred))
28     correct = np.squeeze(correct_tensor.numpy())
29     num_correct += np.sum(correct)
30
31
32 print("Test loss: {:.3f}".format(np.mean(test_losses)))
33 test_acc = num_correct/len(test_loader.dataset)
34 print("Test accuracy: {:.3f}".format(test_acc))

```

Test loss: 0.391
 Test accuracy: 0.844