lstm-sentiment

November 24, 2025

1 Sentiment Classification with RNNs (LSTMs)

In this assignment you will experiment with training and evaluating sentiment classification models that use recurrent neural networks (RNNs) implemented in PyTorch, where an input document (movie review) is represented as a sequence of word embeddings.

If you run the code locally on your computer, you will need to install the PyTorch package, using the instructions shown here.

While knowledge of PyTorch and NumPy is useful, it is not essential for completing this assignment.

1.1 Write Your Name Here:

2 Submission Instructions

While the code in this notebook can be run locally on a powerful machine, it is highly recommended that the notebook is run on the GPU infrastructure available for free through the Educational cluster or in Google's Colab.

2.0.1 Local machine:

- 1. Click the Save button at the top of the Jupyter Notebook.
- 2. Please make sure to have entered your name above.
- 3. Select Cell -> All Output -> Clear. This will clear all the outputs from all cells (but will keep the content of ll cells).
- 4. Select Cell -> Run All. This will run all the cells in order, and will take several minutes.
- 5. Once you've rerun everything, select File -> Save and Export Notebook as -> PDF and download a PDF version *lstm-sentiment.pdf* showing the code and the output of all cells, and save it in the same folder that contains the notebook file *lstm-sentiment.ipynb*.
- 6. Look at the PDF file and make sure all your solutions are there, displayed correctly.
- 7. Submit **both** your PDF and notebook on Canvas. Make sure the PDF and notebook show the outputs of the training and evaluation procedures. Also upload the **output** on the test datasets.
- 8. Verify your Canvas submission contains the correct files by downloading them after posting them on Canvas.

2.0.2 Educational HPC cluster:

1. Once you edited your code using Jupyter notebook, download the Python source code lstm-sentiment.py by selecting File -> Save and Export Notebook as -> Executable Script.

- 2. Run the Python source code on the cluster, using the instructions at https://webpages.charlotte.edu/rbunescu/courses/itcs4101/centaurus.pdf
- 3. Look at the Slurm output file and make sure all your solutions are there, displayed correctly.
- 4. Edit the Analysis section in the notebook file, and save it as a PDF. Alternatively, you can use a text editor to edit your Analysis, then export it as PDF.
- 5. Submit the **Slurm output file**, the **Python source code** file lstm-sentiment.py, the corresponding **Jupyter notebook** file, and the **analysis PDF** on Canvas. Also upload the **output** on the test datasets.
- 6. Verify your Canvas submission contains the correct files by downloading them after posting them on Canvas.

2.0.3 Google Colab:

To load the notebook in Colab:

- 1. Point your browser to https://colab.research.google.com/
- 2. If a pop-up window opens, click on the Upload button and select this notebook file code/lstm-sentiment.ipynb from the homework folder.
- 3. Alternatively, in the notebook window menu clik File -> Open notebook and load the same notebook file.
- 4. You will also need to upload the data folder and the auxiliarry *.py files from the code folder.

 To do this:
 - Using the menu, click 'File' / 'Locate in Drive'. This will open a new browser window showing the contents of the Drive folder containing the notebook file.
 - Using the manu pane on the left, clik on '+ New', followed by 'File upload' (to upload the .py files) or 'Folder upload' (to upload the data folder).
- 5. Select Runtime -> Run all. This will run all the cells in order, and will take several minutes.
- 6. Once you've rerun everything, select File -> Download -> Download .ipynb to save as. a Juupyter notebook. Then select File -> Print -> Save as PDF to save a PDF copy.
- 7. Submit **both** your PDF and notebook on Canvas. Make sure the PDF and notebook show the outputs of the training and evaluation procedures. Also upload the **output** on the test datasets.
- 8. Verify your Canvas submission contains the correct files by downloading them after posting them on Canvas.

3 Vanilla RNN exercise (30p)

Assume a simple one-layer RNN whose computations depends on three parameters w, u, and v as follows: * A linear state update function $h_t = w * h_{t-1} + u * x_t$. * An output function $y_t = \sigma(v * h_t)$.

Find values for the initial state h_0 and the three parameters w, u, and v, such that $y_t \ge 0.5$ if and only if $\sum_{i=1}^t x_t \le 0$.

3.0.1 RNN Solution:

Please use proper Latex or Jupyter-notebook markdown to make your formula's look readable, e.g. subscripts, summation symbols, ...

4 LSTM-based training and evaluation procedures

We will use the RNNet class defined in models.py that uses LSTMs implemented in PyTorch. Depending on the options, this class runs one LSTM (forward) or two LSTMs (bidirectional, forward-backward) on the padded input text. The last state (or concatenated last states), or the average of the states, is used as input to a fully connected network with 3 hidden layers, with a final output sigmoid node computing the probability of the positive class.

First, the necessary imports.

```
[1]: from models import *
    from sentiment_data import *

import random
    import numpy as np
    import torch
    from typing import NamedTuple

class HyperParams(NamedTuple):
        lstm_size: int
        hidden_size: int
        lstm_layers: int
        drop_out: float
        num_epochs: int
        batch_size: int
        seq_max_len: int
```

Here is the training procedure.

```
[2]: # Training procedure for LSTM-based models
     def train_model(hp: HyperParams,
                     train_exs: List[SentimentExample],
                     dev exs: List[SentimentExample],
                     test_exs: List[SentimentExample],
                     word_vectors: WordEmbeddings,
                     use_average, bidirectional):
         train_size = len(train_exs)
         class num = 1
         # Specify training on gpu: set to False to train on cpu
         \# use_gpu = False
         use_gpu = torch.cuda.is_available()
         if use_gpu: # Set tensor type when using GPU
             float_type = torch.cuda.FloatTensor
         else: # Set tensor type when using CPU
             float_type = torch.FloatTensor
         # To get you started off, we'll pad the training input to 60 words to make
      ⇒it a square matrix.
```

```
train_mat = np.asarray([pad_to_length(np.array(ex.indexed_words), hp.
⇒seq_max_len) for ex in train_exs])
   # Also store the actual sequence lengths.
  train_seq_lens = np.array([len(ex.indexed_words) for ex in train_exs])
  # Training input reversed, useful is using bidirectional LSTM.
  train_mat_rev = np.asarray([pad_to_length(np.array(ex.
aget_indexed_words_reversed()), hp.seq_max_len) for ex in train_exs])
  # Extract labels.
  train_labels_arr = np.array([ex.label for ex in train_exs])
  targets = train labels arr
  # Extract embedding vectors.
  embed_size = word_vectors.get_embedding_length()
  embeddings_vec = np.array(word_vectors.vectors).astype(float)
  # Create RNN model.
  rnnModel = RNNet(hp.lstm_size, hp.hidden_size, hp.lstm_layers, hp.drop_out,
                    class_num, word_vectors,
                    use_average, bidirectional,
                    use_gpu =use_gpu)
  # If GPU is available, then run experiments on GPU
  if use_gpu:
      rnnModel.cuda()
  # Specify optimizer.
  optimizer = optim.Adam(filter(lambda p: p.requires_grad, rnnModel.
→parameters()),
                         lr = 5e-3, weight_decay = 5e-3, betas = (0.9, 0.9))
   # Define loss function: Binary Cross Entropy loss for logistic regression
⇔ (binary classification).
  criterion = nn.BCELoss()
   # Get embeddings of words for forward and reverse sentence: (num_ex *
⇒seq_max_len * embedding_size)
  x = np.zeros((train_size, hp.seq_max_len, embed_size))
  x_rev = np.zeros((train_size, hp.seq_max_len, embed_size))
  for i in range(train_size):
      x[i] = embeddings_vec[train_mat[i].astype(int)]
      x_rev[i] = embeddings_vec[train_mat_rev[i].astype(int)]
  # Train the RNN model, gradient descent loop over minibatches.
  for epoch in range(hp.num_epochs):
```

```
rnnModel.train()
      ex_idxs = [i for i in range(train_size)]
      random.shuffle(ex_idxs)
      total_loss = 0.0
      start = 0
      while start < train_size:</pre>
          end = min(start + hp.batch_size, train_size)
           # Get embeddings of words for forward and reverse sentence: (num ex_{f L}
→* seq_max_len * embedding_size)
          x_batch = form_input(x[ex_idxs[start:end]]).type(float_type)
          x_batch_rev = form_input(x_rev[ex_idxs[start:end]]).type(float_type)
          y_batch = form_input(targets[ex_idxs[start:end]]).type(float_type)
          seq_lens_batch = train_seq_lens[ex_idxs[start:end]]
           # Compute output probabilities over all examples in minibatch.
          probs = rnnModel(x_batch, x_batch_rev, seq_lens_batch).flatten()
           # Compute loss over all examples in minibatch.
          loss = criterion(probs, y_batch)
          total_loss += loss.data
           # Zero gradients, perform a backward pass, and update the weights.
          optimizer.zero_grad()
          loss.backward()
          optimizer.step()
          start = end
      print("Loss on epoch %i: %f" % (epoch, total_loss))
      # Print accuracy on training and development data.
      if epoch % 10 == 0:
          acc = eval_model(rnnModel, train_exs, embeddings_vec, hp.
→seq_max_len)
          print('Epoch', epoch, ': Accuracy on training set:', acc)
          acc = eval_model(rnnModel, dev_exs, embeddings_vec, hp.seq_max_len)
          print('Epoch', epoch, ': Accuracy on development set:', acc)
  # Evaluate model on the training dataset.
  acc = eval_model(rnnModel, train_exs, embeddings_vec, hp.seq_max_len)
  print('Accuracy on training set:', acc)
  # Evaluate model on the development dataset.
  acc = eval_model(rnnModel, dev_exs, embeddings_vec, hp.seq_max_len)
```

```
print('Accuracy on develpment set:', acc)
return rnnModel
```

Here is the testing (evaluation) procedure.

```
[3]: # Evaluate the trained model on test examples and return predicted labels on
      ⇔accuracy.
     def eval model(model, exs, embeddings_vec, seq_max_len, pred_only = False):
         # Put model in evaluation mode.
         model.eval()
         # Extract size pf word embedding.
         embed_size = len(embeddings_vec[0])
         # Get embeddings of words for forward and reverse sentence: (num ex *
      ⇒seq_max_len * embedding_size)
         exs_mat = np.asarray([pad_to_length(np.array(ex.indexed_words),__
      ⇒seq_max_len) for ex in exs])
         exs_mat_rev = np.asarray([pad_to_length(np.array(ex.

get_indexed_words_reversed()), seq_max_len) for ex in exs])
         exs_seq_lens = np.array([len(ex.indexed_words) for ex in exs])
         # Get embeddings of words for forward and reverse sentence: (num_ex *
      ⇒seq_max_len * embedding_size)
         x = np.zeros((len(exs), seq_max_len, embed_size))
         x_rev = np.zeros((len(exs), seq_max_len, embed_size))
         for i,ex in enumerate(exs):
             x[i] = embeddings vec[exs mat[i].astype(int)]
             x_rev[i] = embeddings_vec[exs_mat_rev[i].astype(int)]
         x = form input(x)
         x_rev = form_input(x_rev)
         # Run the model on the test examples.
         preds = model(x, x_rev, exs_seq_lens).cpu().detach().numpy().flatten()
         preds[preds >= 0.5] = 1
         preds[preds < 0.5] = 0
         if pred_only == True:
            return preds
         else:
            targets = np.array([ex.label for ex in exs])
            return np.mean(preds == targets)
```

5 Experimental evaluations on the Rotten Tomatoes dataset (15 + 15 + 15p)

First, code for reading the examples and the corresponding GloVe word embeddings.

```
[4]: random.seed(1)
     np.random.seed(1)
     torch.manual seed(1)
     word_vecs_path = '../data/glove.6B.300d-relativized.txt'
     train_path = '../data/rt/train.txt'
     dev_path = '../data/rt/dev.txt'
     blind_test_path = '../data/rt/test-blind.txt'
     test_output_path = 'test-blind.output.txt'
     word_vectors = read_word_embeddings(word_vecs_path)
     word_indexer = word_vectors.word_indexer
     train_exs = read_and_index_sentiment_examples(train_path, word_indexer)
     dev exs = read and index sentiment examples(dev path, word indexer)
     test_exs = read_and_index_sentiment_examples(blind_test_path, word_indexer)
     print(repr(len(train exs)) + " / " +
           repr(len(dev_exs)) + " / " +
           repr(len(test_exs)) + " train / dev / test examples")
```

Read in 30135 vectors of size 300 8530 / 1066 / 1066 train / dev / test examples

5.1 Use only the last state from one LSTM

Evaluate One LSTM + fully connected network, use the last hidden state of LSTM. To get initial results faster, you can try reducing lstm_size, hidden_size, batch_size and even num_epochs.

The accuracy on development data is: *75.61% if trained on my MacBook Pro M1. *77.67% if trained on the HPC educational cluster.

Thus, although using the same random number generator seeds, the actual value may vary depending on machine and version of PyTorch or NumPy.

```
num_epochs = 50, # number of epochs for SGD-based procedure
                 batch_size = 1024, # examples in a minibatch
                  seq_max_len = 60) # maximum length of an example sequence
use_average = False
bidirectional = False
# Train RNN model.
model1 = train_model(hp, train_exs, dev_exs, test_exs, word_vectors,_
 →use_average, bidirectional)
# Generate RNN model predictions for test set.
embeddings_vec = np.array(word_vectors.vectors).astype(float)
test_exs_predicted = eval_model(model1, test_exs, embeddings_vec, hp.
 ⇒seq_max_len, pred_only = True)
# Write the test set output
for i, ex in enumerate(test_exs):
    ex.label = int(test_exs_predicted[i])
write_sentiment_examples(test_exs, test_output_path, word_indexer)
print("Prediction written to file for Rotten Tomatoes dataset.")
Loss on epoch 0: 6.302687
Epoch 0 : Accuracy on training set: 0.551348182883939
Epoch 0: Accuracy on development set: 0.5694183864915572
Loss on epoch 1: 5.987905
Loss on epoch 2: 5.239053
Loss on epoch 3: 4.788215
Loss on epoch 4: 4.641677
Loss on epoch 5: 4.505709
Loss on epoch 6: 4.308260
Loss on epoch 7: 4.272375
Loss on epoch 8: 4.245881
Loss on epoch 9: 4.343611
Loss on epoch 10: 4.275553
Epoch 10 : Accuracy on training set: 0.7875732708089097
Epoch 10: Accuracy on development set: 0.7392120075046904
Loss on epoch 11: 4.097285
Loss on epoch 12: 4.144544
Loss on epoch 13: 4.183259
Loss on epoch 14: 4.015156
Loss on epoch 15: 3.979462
Loss on epoch 16: 4.016921
Loss on epoch 17: 4.036207
Loss on epoch 18: 3.945005
Loss on epoch 19: 3.871135
Loss on epoch 20: 3.913440
```

```
Epoch 20 : Accuracy on training set: 0.7960140679953107
Epoch 20: Accuracy on development set: 0.7420262664165104
Loss on epoch 21: 3.864674
Loss on epoch 22: 3.738769
Loss on epoch 23: 3.718158
Loss on epoch 24: 3.827707
Loss on epoch 25: 3.759531
Loss on epoch 26: 3.722979
Loss on epoch 27: 3.601128
Loss on epoch 28: 3.652118
Loss on epoch 29: 3.553334
Loss on epoch 30: 3.553097
Epoch 30 : Accuracy on training set: 0.8283704572098476
Epoch 30: Accuracy on development set: 0.7617260787992496
Loss on epoch 31: 3.585951
Loss on epoch 32: 3.462646
Loss on epoch 33: 3.517080
Loss on epoch 34: 3.494456
Loss on epoch 35: 3.453969
Loss on epoch 36: 3.475022
Loss on epoch 37: 3.433378
Loss on epoch 38: 3.296738
Loss on epoch 39: 3.396697
Loss on epoch 40: 3.401010
Epoch 40: Accuracy on training set: 0.8355216881594373
Epoch 40: Accuracy on development set: 0.7607879924953096
Loss on epoch 41: 3.393050
Loss on epoch 42: 3.328170
Loss on epoch 43: 3.347071
Loss on epoch 44: 3.207592
Loss on epoch 45: 3.419304
Loss on epoch 46: 3.187574
Loss on epoch 47: 3.089619
Loss on epoch 48: 3.125739
Loss on epoch 49: 3.190348
Accuracy on training set: 0.8576787807737397
Accuracy on develpment set: 0.7560975609756098
Prediction written to file for Rotten Tomatoes dataset.
```

5.2 Use the average of all states from one LSTM (15p)

Evaluate One LSTM + fully connected network, use average of all states of the LSTM.

Our accuracy on development data is 77.67%

```
[]: random.seed(1)
   np.random.seed(1)
   torch.manual_seed(1)
```

```
## YOUR CODE HERE
```

5.3 Use a bidirectional LSTM, concatenate last states (15p)

Evaluate Two LSTMs (bidirectional) + fully connected network, concatenate their last states.

Our accuracy on development data is 76.83%

```
[]: random.seed(1)
    np.random.seed(1)
    torch.manual_seed(1)

## YOUR CODE HERE
```

5.4 Use a bidirectional LSTM, concatenate the averages of their states (15p)

Evaluate Two LSTMs (bidirectional) + fully connected network, concatenate the averages of their states.

Our accuracy on development data is 77.39%

```
[]: random.seed(1)
    np.random.seed(1)
    torch.manual_seed(1)

## YOUR CODE HERE
```

5.5 [Bonus] Average performance and standard deviation (20p)

The NN performance can vary depending on the random initialization of its parameters. Train and evaluate each model 10 times, from different random initializations (10 different seeds). Average the accuracy over the 10 runs and compare the performance of the 4 models on the Rotten Tomatoes dataset. Report in your analysis the average and standard deviation for each model.

```
[ ]: ## YOUR CODE HERE
```

6 Experimental evaluations on the IMDB dataset (10 + 10 + 10 + 10p)

Run the same 4 evaluations on the IMDB dataset (10p for each evaluation).

```
[]: train_path = '../data/imdb/train.txt'
  dev_path = '../data/imdb/dev.txt'
  test_path = '../data/imdb/test.txt'

test_output_path = 'test-imdb.output.txt'

## YOUR CODE HERE
```

6.1 [Bonus] Cross-domain performance (15p)

Compare the performance of the Bidirectional LSTM with state averaging on the IMDB test set in two scenarios:

- 1. The model is trained on the IMDB training data.
- 2. The model is trained on the Rotten Tomatoes data.

[]: ## YOUR CODE HERE

6.2 [Bonus] Base rate neglect cognitive bias (15p)

Read the lecture slides on Bias, Fairness, and Rationality and solve the problem on slide 12 on the base rate neglect cognitive bias. For your convenience, the problem is copied below.

Context: 10 out of every 1,000 women at age forty who participate in routine screen have breast cancer. 8 of every 10 women with breast cancer will get a positive mammography. 95 out of every 990 women without breast cancer will also get a positive mammography.

Question: Consider a new representative sample of women at age forty who got a positive mammography in routine screening. What is their likelihood of having breast cancer?

6.2.1 Your Solution Here:

Please use proper Latex or Jupyter-notebook markdown to make your formula's look readable, e.g. subscripts, summation symbols, ...

6.3 Bonus points

Anything extra goes here.

6.4 Analysis (20p)

Include an analysis of the results that you obtained in the experiments above. Also compare with the sentiment classification performance from previous assignments and explain the difference in accuracy. Show the results using table(s).