# Uppsala University

## Deep learning

---

# Hand-in assignment (3)

---

Tong You

June 10, 2021

# Contents

# 1  Introduction

Whereas in the previous two assignments we were working with image data (MNIST and Warwick data), in this final assignment we will implement a language model. We will be using the PennTreebank (PTB) as our dataset. This is a collection of articles from the Wall Street Journal. We will implement the language model in several stages. First, we need to preprocess the data. Next we will implement a (simple) Elmann RNN. Finally, we will look at several different ways to improve the model. We will show generated example phrases at each stage in order to show the differences. The model will be implemented in Pytorch using an RTX 2080 Ti in order to do GPU training.

# 2  The dataset

As mentioned previously, we will be using the PTB dataset. This is split up into three splits: a training split, validation split, and testing split. Before we can feed this data into our language model we need to preprocess the data, generate a dictionary with our keys being the unique words in the three splits and the values being an unique integer corresponding to a particular unique word. Before we proceed, we will be showing some example sentences. We pick two random sentences from each split.

- Example training set phrase: in addition the cray-3 will contain N processors twice as many as the largest current supercomputer

- Example training set phrase: what 's more the test and learning materials are both produced by the same company <unk> a joint venture of mcgraw-hill inc. and macmillan 's parent britain 's maxwell communication corp

- Example validation set phrase: the ghost of the soviet <unk> discovered in cuba back in the <unk> costs just a few hundred million the price of the caribbean command in key west that president carter created in N

- Example validation set phrase: mr. wolf has <unk> merger advice from a major wall street securities firm relying instead only on a takeover lawyer peter <unk> of <unk> <unk> slate <unk> & flom

- Example test set phrase: the executive denied speculation that saatchi was bringing in the new chief executive officer only to clean up the company financially so that the brothers could lead a buy-back

- Example test set phrase: quantum 's lot is mostly tied to polyethylene <unk> used to make garbage bags milk <unk> <unk> toys and meat packaging among other items

The next step is to add an additional term <eos> that denotes the end of a sentence. All code will be attached in the Appendix. Next we want to check the size of the three splits. For this count we consider the special characters <eos> and <unk>. The size of the splits:

- Number of training words: 929589

- Number of validation words: 73760

- Number of testing words: 82430

From the above we want to generate a dictionary as mentioned earlier. In order to generate the dictionary we first flatten the three splits and combine the splits into one big list. Then we use the set function in Python. This will generate a set with the unique words. From this we can easily see how many unique words there are. There are 10000 unique words in the entire PTB dataset. Finally, we use this dictionary to map our words to a unique integer. To generate sample phrases we would look up the word corresponding to any unique integer. Our input to the RNN will therefore consist of a vector of integers. Next we will look at the language model.

# 3 A simple language model

In the previous section we did some preprocessing on the PTB dataset. Now we are ready to implement a simple RNN using PyTorch. For implementation details we refer to the Appendix, however code that is deemed important is explained in the following text. After some initial tweaking, an "ideal" set of hyperparameters was chosen. The test perplexity still doesn't reach the value of a highly-optimized Elman RNN. However, we will implement some methods to try to lower the perplexity. In order to make an one-to-one comparison between the sections we will use the same input sentence(s) when generating sample phrases.

## 3.1 Results

To start with, we will show results of a vanilla RNN without any gradient clipping. In order to highlight differences between different methods we will try to use the same parameters in order to really show the gains that were made. With normal SGD it was very difficult get a low perplexity that starts to approach a well-tuned Elmann RNN. The following loss curves correspond to a training batch size of 32 and validation/testing batch sizes of 512. The sequence length for all splits was 50. Both the hidden and embedding dimensions were set to 500. The RNN was implemented with a layer number of 1 first. The loss was averaged each epoch to yield less noisier training/validation curves. We will also use dropout since that really improved the training for the simple RNN (value of 0.5). To start with, we will show results of a vanilla RNN without any gradient clipping. In order to highlight differences between different methods we will try to use the same parameters in order to really show the gains that were made. With normal SGD it was very difficult get a low perplexity that starts to approach a well-tuned Elmann RNN. The following loss curves correspond to a training batch size of 32 and validation/testing batch sizes of 512. The sequence length for all splits was 50. Both the hidden and embedding dimensions were set to 500. The RNN was implemented with a layer number of 1 first. The loss was averaged each epoch to yield less noisier training/validation curves. We will also use dropout since that really improved the training for the simple RNN (value of 0.5). The learning rate was 0.5. The loss curve below shows the results for the simple RNN. This model was then evaluated on the test split and we got a perplexity of 214.15121520493514. This is still quite high, but not as high as initial testing of parameters.
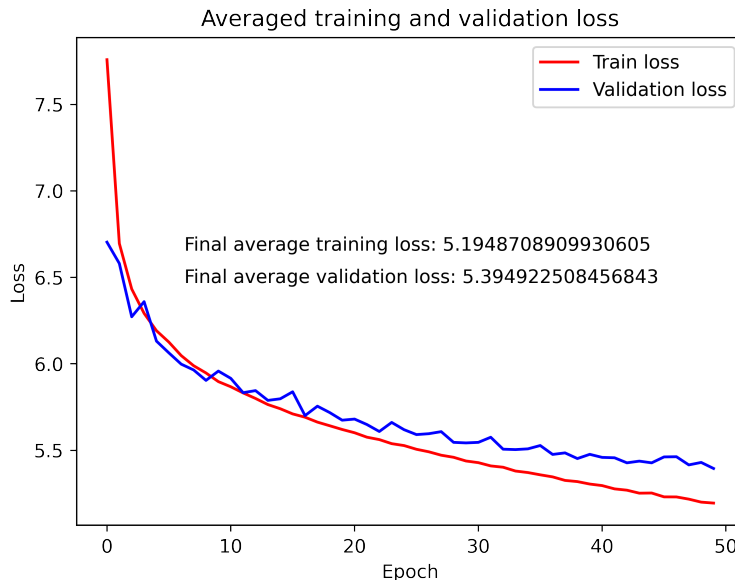
Figure 1: Training and validation loss averaged over all iterations each epoch. The final losses are annotated in the figure. The RNN has one layer.

We will see in the next section how we can improve this.

# 4 Improving the recurrent model

As mentioned earlier, there are several ways to improve upon the previous result. In addition we will also show examples of generated phrases by providing a seeding input to the trained model.

## 4.1 Results - gradient clipping

We will first show some results of gradient clipping applied to a two-layer RNN. We will use the same parameters, but after some preliminary testing we found that we can use quite high learning rates, so we chose 10. We clipped the gradients at 0.5 using the L-2 norm. The one-layer RNN got higher perplexity than the one-layer network with perplexity in testing so we decided to show results for a two-layer network with gradient clipping.

This network gets higher test perplexity than the above network. The test perplexity after 50 epochs was 247.3608380295578.
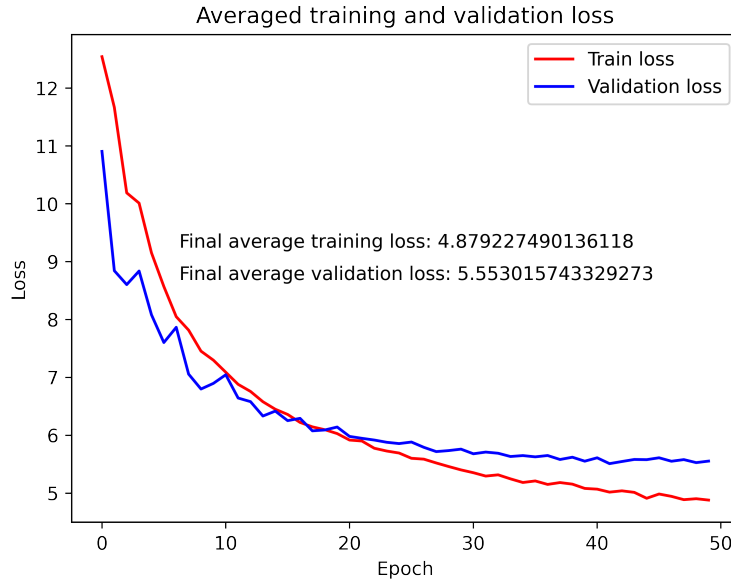


Figure 2: Training and validation loss averaged over all iterations each epoch. Training was done with gradient clipping. The final losses are annotated in the figure. The RNN has two layers.

Next we will show a generated phrase giving an input sequence. We will use a greedy method. However, the current model outputs the same word when we specify how many words we want to generate. Several attempts to fix this didn't work, so only one generated phrase is shown. However, it seems that the model outputs the most frequent word. The seeding sequence is 2 words long and is created by using uniformly generated random integers up to the number of unique words. We started with "there" and "unrelated".

- Generated phrase: "there" "unrelated" "with" "with" "with"

## 4.2 Results - LSTM

Since the Elmann RNN is quite "simple", we will next use an LSTM. These are much less prone to exploding gradients, although we didn't encounter

6

such problems for our implementation. We will use the exact same parameters as the results in the subsection above, only changing the RNN to an LSTM. For ease of reading, we will give them again. For the training we used a batch size of 32 and for validation/testing we evaluated on data with batch sizes of 512. The sequences were of length 50 for all splits. We also used the same embedding and hidden dimensions of 500. And the LSTM in this case has two layers. We also found that a dropout of 0.8 yielded quite good results. The learning rate was set to 10.0 with gradient clipping with a max L2-norm of 0.5. We also trained for 50 epochs using standard SGD. The final test perplexity was 159.23738342645728.
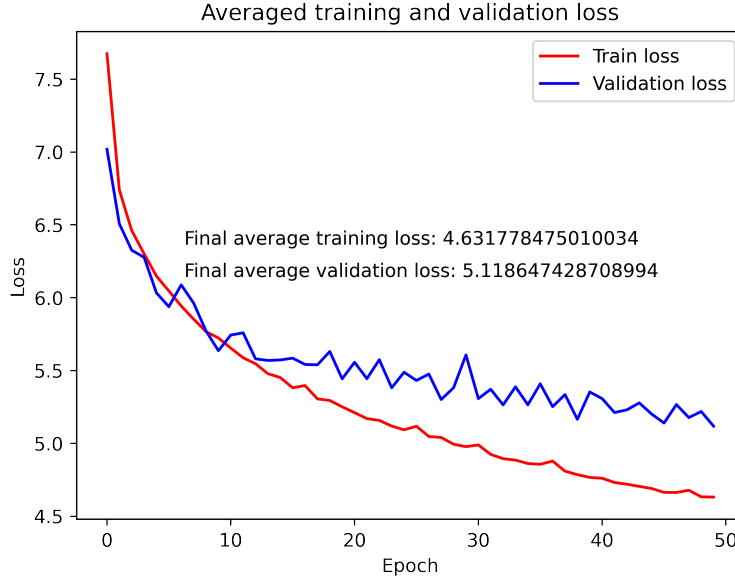


Figure 3: Training and validation loss averaged over all iterations each epoch. Training was done with gradient clipping. The final losses are annotated in the figure. Instead of a two-layer RNN, a two-layer LSTM was trained.

Next we show an example of a generated phrase. We used as our input sequence: "there" and "competitiveness". We wanted to generate three words.

- Generated phrase: "there" "competitiveness" "is" "<unk>" "<unk>".

## 4.3 Results - additional modifications

There are different changes we can apply to our current model and see if we can spot any improvements. We will use the same settings as used above, but swap the LSTM with an GRU. Training will be done in the same way.
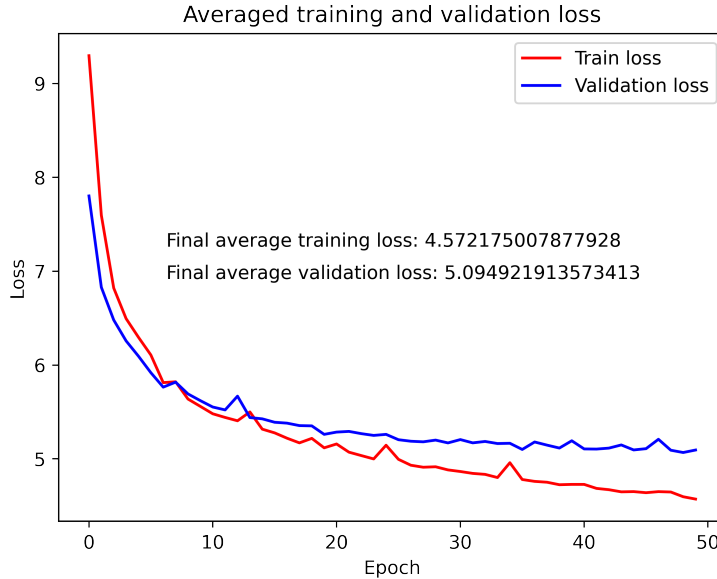


Figure 4: Training and validation loss averaged over all iterations each epoch. Training was done with gradient clipping. The final losses are annotated in the figure. A two-layer GRU was used to train the PTB dataset.

The final test perplexity was slightly lower than the LSTM: 157.50157610436733. An example of a generated phrase:

- Generated phrase: "there" "unknown" "is" "a" "<unk>".

The trained model was seeded with a sequence in a similar way as above. The model generates much more promising results, however it still outputs the same word after a few generations. We will finally try to improve on the previous model by using a custom learning rate scheduler. We will still use the GRU but instead use a cosine learning rate scheduler and train for 50 epochs again. We will use SGD again with a learning rate of 10, but this

will be modified by the cosine learning rate scheduler. The scheduler was implemented as follows:

```
cosine_sched = optim.lr_scheduler.CosineAnnealingLR(optimizer,
↪   T_max = 500, eta_min = 0.1, last_epoch = -1, verbose =
↪   True)
```
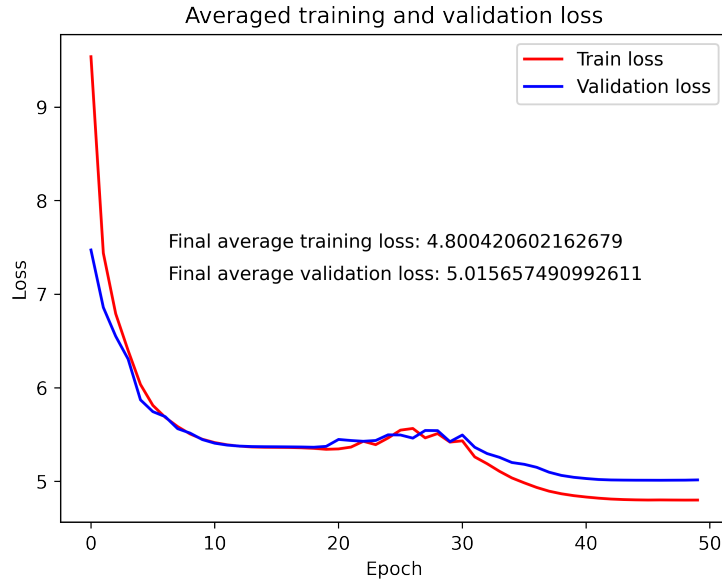
Figure 5: Training and validation loss averaged over all iterations each epoch. Training was done with gradient clipping. The final losses are annotated in the figure. A two-layer GRU was used to train the PTB dataset. This time a cosine learning rate scheduler was used.

With the same parameters as the GRU with SGD and no learning rate scheduling, we now get a final perplexity of 146.02419039272354. This is about an 11 point difference compared to the GRU without rate scheduling. We can see in the loss curves that the loss peaked a bit after about 25 epochs, but decreased again afterwards. Better results might be expected by combining an exponential decay scheduler with a cosine learning rate scheduler. Despite the improvements in perplexity using the same seeding mechanism

as we had previously the generated phrases still repeat the same word after only a few words:

- Generated phrase: "there" "brewing" "is" "a" "<unk>".

# 5 Appendix

Since the code was implemented in a single Jupyter Notebook, the code in the different cells will be attached here as a single file.

```python
import torch
import torch.nn as nn
import torch.optim as optim
import torch.nn.functional as F

import functools
import operator
import random
import pickle

import numpy as np

from sklearn.metrics import confusion_matrix
from torch.utils.data import TensorDataset
from torch.utils.data import DataLoader
from matplotlib import pyplot as plt

# Exercise 1
# Training data
keepNew = True
with open('PTB/ptb.train.txt', mode = 'r', newline = None) as
    train_f:
    train_dat = train_f.read().splitlines(keepNew)

# Validation data
with open('PTB/ptb.valid.txt', mode = 'r', newline = None) as
    valid_f:
    valid_dat = valid_f.read().splitlines(keepNew)
```

```python
27
28  # Testing data
29  with open('PTB/ptb.test.txt', mode = 'r', newline = None) as
    ↪   test_f:
30      test_dat = test_f.read().splitlines(keepNew)
31
32  # Printing examples phrases
33  n_examp = np.random.choice(1000, size = 6, replace = False)
34
35  print("Example training set phrase:", train_dat[n_examp[0]])
36  print("Example training set phrase:", train_dat[n_examp[1]])
37  print("Example validation set phrase:", valid_dat[n_examp[2]])
38  print("Example validation set phrase:", valid_dat[n_examp[3]])
39  print("Example test set phrase:", test_dat[n_examp[4]])
40  print("Example test set phrase:", test_dat[n_examp[5]])
41
42  # Adding <eos>
43  train_proc, valid_proc, test_proc = [], [], []
44  for l_train in train_dat:
45      l_mod = l_train.replace('\n', '<eos>')
46      train_proc.append(l_mod)
47
48  for l_valid in valid_dat:
49      l_mod = l_valid.replace('\n', '<eos>')
50      valid_proc.append(l_mod)
51
52  for l_train in test_dat:
53      l_mod = l_train.replace('\n', '<eos>')
54      test_proc.append(l_mod)
55
56  # Split up each line in individual words
57  train_words, valid_words, test_words = [], [], []
58  for tp in train_proc:
59      train_words.append(tp.split())
60
61  for tp in valid_proc:
62      valid_words.append(tp.split())
63
```

```
64  for tp in test_proc:
65      test_words.append(tp.split())
66
67  # Flatten list of lists into a single list
68  train_words = functools.reduce(operator.iconcat, train_words,
    ↪   [])
69  valid_words = functools.reduce(operator.iconcat, valid_words,
    ↪   [])
70  test_words = functools.reduce(operator.iconcat, test_words,
    ↪   [])
71
72  num_train = len(train_words)
73  num_valid = len(valid_words)
74  num_test = len(test_words)
75
76  print("Number of training words:", num_train)
77  print("Number of validation words:", num_valid)
78  print("Number of testing words:", num_test)
79
80  # Building a dictionary
81  all_words = train_words + valid_words + test_words
82  set_words = set(all_words)
83  num_unique = len(set_words)
84
85  print("Number of unique words in training + validation +
    ↪   testing splits:", num_unique)
86
87  num_id = np.random.choice(num_unique, size = num_unique,
    ↪   replace = False)
88
89  n = 0
90  unique_dict = {}
91  for uw in set_words:
92      unique_dict.update({uw : num_id[n]})
93      n += 1
94
95  # Replacing all words in the training/validation/testing
    ↪   splits with their integer representation
```

```python
96  train_ints, valid_ints, test_ints = [], [], []
97
98  for word in train_words:
99      int_rep = unique_dict[word]
100     train_ints.append(int_rep)
101
102 for word in valid_words:
103     int_rep = unique_dict[word]
104     valid_ints.append(int_rep)
105
106 for word in test_words:
107     int_rep = unique_dict[word]
108     test_ints.append(int_rep)
109
110 # Check if CUDA is available
111 device = torch.device("cuda") if torch.cuda.is_available()
    ↪   else torch.device("cpu")
112 print("Device:", torch.cuda.get_device_name(device))
113
114 # Resetting model function
115 # Credits:
    ↪   https://discuss.pytorch.org/t/reset-model-weights/19180/4
116 def reset_model(model):
117     for layer in model.children():
118         if hasattr(layer, 'reset_parameters'):
119             layer.reset_parameters()
120
121 # Exercise 2
122 # Convert our training/validation/testing splits to Torch
    ↪   tensors
123 train_dat = torch.tensor(train_ints)
124 valid_dat = torch.tensor(valid_ints)
125 test_dat = torch.tensor(test_ints)
126
127 ### Data loading
128 batch_size = 32
129 batch_eval = 512
130 batch_test = 512
```

```
131  seq_train = 50
132  seq_valid = 50
133  seq_test = 50
134
135  s_train_l = num_train // seq_train
136  s_valid_l = num_valid // seq_valid
137  s_test_l = num_test // seq_test
138
139  # Trim training/validation/testing data and reshape into
     ↪   tensor of num_sequences by sequence_length
140  # Training data
141  train_seq = torch.narrow(train_dat, 0, 0, seq_train *
     ↪   s_train_l)
142  train_lab = torch.roll(train_seq, shifts = -1, dims = 0)
143
144  train_seq = train_seq.reshape(s_train_l, seq_train)
145  train_lab_seq = train_lab.reshape(s_train_l, seq_train)
146
147  # Validation data
148  valid_seq = torch.narrow(valid_dat, 0, 0, seq_valid *
     ↪   s_valid_l)
149  valid_lab = torch.roll(valid_seq, shifts = -1, dims = 0)
150
151  valid_seq = valid_seq.reshape(s_valid_l, seq_valid)
152  valid_lab_seq = valid_lab.reshape(s_valid_l, seq_valid)
153
154  # Testing data
155  test_seq = torch.narrow(test_dat, 0, 0, seq_test * s_test_l)
156  test_lab = torch.roll(test_seq, shifts = -1, dims = 0)
157
158  test_seq = test_seq.reshape(s_test_l, seq_test)
159  test_lab_seq = test_lab.reshape(s_test_l, seq_test)
160
161  # Divide training and validation data into correct
     ↪   mini-batches
162  num_batches = train_seq.shape[0] // batch_size
163  valid_batches = valid_seq.shape[0] // batch_eval
164  test_batches = test_seq.shape[0] // batch_test
```

```
165
166  training_set = TensorDataset(train_seq.to(device),
     ↪   train_lab_seq.to(device))
167  training_loader = DataLoader(training_set, shuffle = False,
     ↪   batch_size = num_batches)
168
169  valid_set = TensorDataset(valid_seq.to(device),
     ↪   valid_lab_seq.to(device))
170  valid_loader = DataLoader(valid_set, shuffle = False,
     ↪   batch_size = valid_batches)
171
172  test_set = TensorDataset(test_seq.to(device),
     ↪   test_lab_seq.to(device))
173  test_loader = DataLoader(test_set, shuffle = False, batch_size
     ↪   = test_batches)
174
175  ### RNN code
176  # Embedding parameter
177  embed_dim = 500
178
179  # RNN parameters
180  hidden_dim = 500
181  in_size = embed_dim
182  n_layers = 2
183
184  # Vanilla RNN
185  class ElmanRNN(nn.Module):
186      def __init__(self, input_size, hidden_size, num_layers,
         ↪   num_embeddings, embedding_dim, hidden_dim, num_unique,
         ↪   drop_out):
187          super(ElmanRNN, self).__init__()
188
189          self.input_size = input_size
190          self.hidden_size = hidden_size
191          self.num_layer = num_layers
192          self.drop_out = drop_out
193          self.num_embeddings = num_embeddings
194          self.embedding_dim = embedding_dim
```

15

```python
            self.embed = nn.Embedding(num_embeddings,
            ↪   embedding_dim)
            self.elman = nn.RNN(input_size, hidden_size,
            ↪   num_layers, dropout = drop_out, batch_first =
            ↪   True)
            self.linear = nn.Linear(hidden_dim, num_unique)

    def forward(self, mod_input):
            word_embed = self.embed(mod_input)
            rnn_out, hidden_out = self.elman(word_embed)
            rnn_out = self.linear(rnn_out)
            rnn_out = rnn_out.view(-1, num_unique)

            return rnn_out, hidden_out

# LSTM
class LSTMRNN(nn.Module):
    def __init__(self, input_size, hidden_size, num_layers,
    ↪   num_embeddings, embedding_dim, hidden_dim, num_unique,
    ↪   drop_out):
            super(LSTMRNN, self).__init__()

            self.input_size = input_size
            self.hidden_size = hidden_size
            self.num_layer = num_layers
            self.drop_out = drop_out
            self.num_embeddings = num_embeddings
            self.embedding_dim = embedding_dim

            self.embed = nn.Embedding(num_embeddings,
            ↪   embedding_dim)
            self.lstm = nn.LSTM(input_size, hidden_size,
            ↪   num_layers, dropout = drop_out, batch_first =
            ↪   True, bidirectional = False)
            self.linear = nn.Linear(hidden_dim, num_unique)

    def forward(self, mod_input):
```

```
225         word_embed = self.embed(mod_input)
226         lstm_out, hidden_out = self.lstm(word_embed)
227         lstm_out = self.linear(lstm_out)
228         lstm_out = lstm_out.view(-1, num_unique)
229
230         return lstm_out, hidden_out
231
232 # GRU
233 class GRURNN(nn.Module):
234     def __init__(self, input_size, hidden_size, num_layers,
        ↪  num_embeddings, embedding_dim, hidden_dim, num_unique,
        ↪  drop_out):
235         super(GRURNN, self).__init__()
236
237         self.input_size = input_size
238         self.hidden_size = hidden_size
239         self.num_layer = num_layers
240         self.drop_out = drop_out
241         self.num_embeddings = num_embeddings
242         self.embedding_dim = embedding_dim
243
244         self.embed = nn.Embedding(num_embeddings,
            ↪  embedding_dim)
245         self.gru = nn.GRU(input_size, hidden_size, num_layers,
            ↪  dropout = drop_out, batch_first = True,
            ↪  bidirectional = False)
246         self.linear = nn.Linear(hidden_dim, num_unique)
247
248     def forward(self, mod_input):
249         word_embed = self.embed(mod_input)
250         gru_out, hidden_out = self.gru(word_embed)
251         gru_out = self.linear(gru_out)
252         gru_out = gru_out.view(-1, num_unique)
253
254         return gru_out, hidden_out
255
256 elman_rnn = ElmanRNN(in_size, hidden_dim, n_layers,
    ↪  num_unique, embed_dim, hidden_dim, num_unique, 0.5)
```

```python
257  lstm_rnn = LSTMRNN(in_size, hidden_dim, n_layers, num_unique,
     ↪   embed_dim, hidden_dim, num_unique, 0.8)
258  gru_rnn = GRURNN(in_size, hidden_dim, n_layers, num_unique,
     ↪   embed_dim, hidden_dim, num_unique, 0.8)
259
260  # Selecting model
261  model = gru_rnn
262
263  reset_model(model)
264  model.to(device)
265
266  # Optimizer
267  l_rate = 10.0
268  sgd = optim.SGD(model.parameters(), lr = l_rate, weight_decay
     ↪   = 0, momentum = 0.0)
269  adam = optim.Adam(model.parameters(), lr = l_rate, betas =
     ↪   (0.9, 0.999), eps = 1e-08, weight_decay = 0.0, amsgrad =
     ↪   False)
270
271  optimizer = sgd
272
273  # Cross entropy loss
274  loss = nn.CrossEntropyLoss()
275
276  # Training parameters
277  num_epochs = 50
278  train_epoch = np.zeros(num_epochs)
279  valid_epoch = np.zeros(num_epochs)
280
281  train_loss = []
282  valid_loss = []
283  test_loss = []
284
285  # Gradient clipping
286  clipGrad = True
287
288  # Learning rate scheduler
```

```python
289  exp_sched = optim.lr_scheduler.ExponentialLR(optimizer, gamma
     ↪   = 0.5,last_epoch = -1, verbose = True)
290  cosine_sched = optim.lr_scheduler.CosineAnnealingLR(optimizer,
     ↪   T_max = 500, eta_min = 0.1, last_epoch = -1, verbose =
     ↪   True)
291  #plat_sched =
     ↪   optim.lr_scheduler.ReduceLROnPlateau(optimizer, mode =
     ↪   'min', factor = 0.1, patience = 10, threshold = 0.0001,
     ↪   threshold_mode = 'rel', cooldown = 0, min_lr = 0, eps =
     ↪   1e-08, verbose = True)
292
293  sched = cosine_sched
294  useScheduler = True
295
296  # Training and evaluation on validation data
297  for epoch in range(num_epochs):
298      train_losss = []
299      valid_losss = []
300
301      for data, lab_train in training_loader:
302          model.train()
303          optimizer.zero_grad()
304
305          pred_out, hid_out = model(data)
306
307          ce_loss = loss(pred_out, lab_train.flatten())
308
309          train_loss.append(ce_loss.item())
310          train_losss.append(ce_loss.item())
311
312          ce_loss.backward()
313
314          # Using gradient clipping
315          if clipGrad:
316              nn.utils.clip_grad_norm_(model.parameters(),
                 ↪   max_norm = 0.5, norm_type = 2.0)
317
318          optimizer.step()
```

```python
319
320          # Using learning rate scheduler
321          if useScheduler:
322              sched.step()
323
324      train_epoch[epoch] = np.mean(train_losss)
325
326      for data_eval, lab_eval in valid_loader:
327          model.eval()
328          with torch.no_grad():
329              valid_preds, valid_hid = model(data_eval)
330
331              ce_valid = loss(valid_preds, lab_eval.flatten())
332              valid_loss.append(ce_valid.item())
333              valid_losss.append(ce_valid.item())
334
335      valid_epoch[epoch] = np.mean(valid_losss)
336
337      print("Epoch: %s" % (epoch + 1))
338
339 # Evaluation on test data
340 model.eval()
341 for data_test, lab_test in test_loader:
342     with torch.no_grad():
343              test_preds, test_hid = model(data_test)
344
345              ce_test = loss(test_preds, lab_test.flatten())
346              test_loss.append(ce_test.item())
347
348 plt.figure(1, figsize = (6.4, 4.8))
349 train, = plt.plot(train_epoch, 'r')
350 valid, = plt.plot(valid_epoch, 'b')
351 plt.xlabel("Epoch")
352 plt.ylabel("Loss")
353 plt.title("Averaged training and validation loss")
354 plt.legend([train, valid], ['Train loss', 'Validation loss'])
```

```python
355  plt.annotate("Final average training loss: %s" %
     ↪  (train_epoch[-1]) ,xycoords = 'figure fraction', xy =
     ↪  (0.25,0.55))
356  plt.annotate("Final average validation loss: %s" %
     ↪  (valid_epoch[-1]), xycoords = 'figure fraction', xy =
     ↪  (0.25,0.50))
357  plt.savefig("train_valid_epoch", dpi = 500)
358  print("Final average training loss: %s." % train_epoch[-1])
359  print("Final average validation loss: %s." % valid_epoch[-1])
360  print()
361
362  plt.figure(2, figsize = (6.4, 4.8))
363  train_it, = plt.plot(train_loss, 'r')
364  valid_it, = plt.plot(valid_loss, 'b')
365  plt.xlabel("Iteration")
366  plt.ylabel("Loss")
367  plt.title("Training and validation loss")
368  plt.legend([train_it, valid_it], ['Train loss', 'Validation
     ↪  loss'])
369  plt.annotate("Final training loss: %s" % (train_loss[-1])
     ↪  ,xycoords = 'figure fraction', xy = (0.25,0.55))
370  plt.annotate("Final validation loss: %s" % (valid_loss[-1]),
     ↪  xycoords = 'figure fraction', xy = (0.25,0.50))
371  #plt.savefig("train_valid_loss_iter", dpi = 500)
372  print("Final training loss: %s." % train_loss[-1])
373  print("Final validation loss: %s." % valid_loss[-1])
374  print()
375
376  plt.figure(3, figsize = (6.4, 4.8))
377  train_fin_ep, = plt.plot(train_losss, 'r')
378  valid_fin_ep, = plt.plot(valid_losss, 'b')
379  plt.xlabel("Iteration")
380  plt.ylabel("Loss")
381  plt.title("Training and validation loss final epoch")
382  plt.legend([train_fin_ep, valid_fin_ep], ['Train loss',
     ↪  'Validation loss'])
```

```python
383 plt.annotate("Training loss final epoch: %s" %
    ↪   (train_losss[-1]) ,xycoords = 'figure fraction', xy =
    ↪   (0.25,0.55))
384 plt.annotate("Validation loss final epoch: %s" %
    ↪   (valid_losss[-1]), xycoords = 'figure fraction', xy =
    ↪   (0.25,0.50))
385 #plt.savefig("train_valid_final_epoch", dpi = 500)
386
387 perplexity = np.exp(np.mean(test_loss))
388 print("Test perplexity: %s." % (perplexity))
389
390 # Generate sample phrases
391 dictionary = unique_dict # number of unique words/tokens is
    ↪   10000
392 dict_keys = list(dictionary.keys())
393 dict_values = list(dictionary.values())
394
395 # Words
396 w_1 = "there"
397 i_1 = dictionary[w_1]
398 w_2 = "is"
399 i_2 = dictionary[w_2]
400 w_3 = "a"
401 i_3 = dictionary[w_3]
402 w_4 = "why"
403 i_4 = dictionary[w_4]
404 w_5 = "on"
405 i_5 = dictionary[w_5]
406
407 # Input sequence
408 input_tens = torch.randint(num_unique, (2,
    ↪   1)).long().to(device)
409 input_tens[0] = i_1
410 print("Input", input_tens)
411
412 inputs_list = input_tens.tolist()
413 print("List input", inputs_list)
414
```

```python
softmax_layer = nn.Softmax(dim = 1)
num_words = 5 # how many words to predict
model.eval()
with torch.no_grad():
    for n in range(num_words):
        out, hidden = model(input_tens)
        probs = softmax_layer(out)

        maxv, maxi = torch.max(probs, dim = 1)
        sortedv, sortedi = -np.sort(-maxv.cpu().numpy()),
            ↪  -np.sort(-maxi.cpu().numpy())

        next_word = sortedi[0]

        inputs_list.append([next_word])

        input_tens = input_tens.tolist()
        input_tens.append([next_word])
        input_tens =
            ↪  torch.tensor(input_tens).long().to(device)


print(inputs_list)
output_sentence = []

for i in inputs_list:
    word_i = dict_keys[dict_values.index(i[0])]
    output_sentence.append(word_i)

print(output_sentence)
```