

1 rnn.py

```
1 from builtins import range
2 from builtins import object
3 import numpy as np
4
5 from cs231n.layers import *
6 from cs231n.rnn_layers import *
7
8
9 class CaptioningRNN(object):
10     """
11     A CaptioningRNN produces captions from image features using a recurrent
12     neural network.
13
14     The RNN receives input vectors of size D, has a vocab size of V, works on
15     sequences of length T, has an RNN hidden dimension of H, uses word vectors
16     of dimension W, and operates on minibatches of size N.
17
18     Note that we don't use any regularization for the CaptioningRNN.
19     """
20
21     def __init__(self, word_to_idx, input_dim=512, wordvec_dim=128,
22                 hidden_dim=128, cell_type='rnn', dtype=np.float32):
23         """
24         Construct a new CaptioningRNN instance.
25
26         Inputs:
27         - word_to_idx: A dictionary giving the vocabulary. It contains V entries,
28           and maps each string to a unique integer in the range [0, V).
29         - input_dim: Dimension D of input image feature vectors.
30         - wordvec_dim: Dimension W of word vectors.
31         - hidden_dim: Dimension H for the hidden state of the RNN.
32         - cell_type: What type of RNN to use; either 'rnn' or 'lstm'.
33         - dtype: numpy datatype to use; use float32 for training and float64 for
34           numeric gradient checking.
35         """
36         if cell_type not in {'rnn', 'lstm'}:
37             raise ValueError('Invalid cell_type "%s"' % cell_type)
38
39         self.cell_type = cell_type
40         self.dtype = dtype
41         self.word_to_idx = word_to_idx
42         self.idx_to_word = {i: w for w, i in word_to_idx.items()}
43         self.params = {}
44
45         vocab_size = len(word_to_idx)
46
47         self._null = word_to_idx['<NULL>']
48         self._start = word_to_idx.get('<START>', None)
49         self._end = word_to_idx.get('<END>', None)
50
51         # Initialize word vectors
52         self.params['W_embed'] = np.random.randn(vocab_size, wordvec_dim)
53         self.params['W_embed'] /= 100
54
55         # Initialize CNN -> hidden state projection parameters
56         self.params['W_proj'] = np.random.randn(input_dim, hidden_dim)
57         self.params['W_proj'] /= np.sqrt(input_dim)
58         self.params['b_proj'] = np.zeros(hidden_dim)
59
60         # Initialize parameters for the RNN
61         dim_mul = {'lstm': 4, 'rnn': 1}[cell_type]
62         self.params['Wx'] = np.random.randn(wordvec_dim, dim_mul * hidden_dim)
63         self.params['Wx'] /= np.sqrt(wordvec_dim)
64         self.params['Wh'] = np.random.randn(hidden_dim, dim_mul * hidden_dim)
65         self.params['Wh'] /= np.sqrt(hidden_dim)
66         self.params['b'] = np.zeros(dim_mul * hidden_dim)
67
68         # Initialize output to vocab weights
69         self.params['W_vocab'] = np.random.randn(hidden_dim, vocab_size)
70         self.params['W_vocab'] /= np.sqrt(hidden_dim)
71         self.params['b_vocab'] = np.zeros(vocab_size)
72
73         # Cast parameters to correct dtype
```

```

74     for k, v in self.params.items():
75         self.params[k] = v.astype(self.dtype)
76
77
78 def loss(self, features, captions):
79     """
80     Compute training-time loss for the RNN. We input image features and
81     ground-truth captions for those images, and use an RNN (or LSTM) to compute
82     loss and gradients on all parameters.
83
84     Inputs:
85     - features: Input image features, of shape (N, D)
86     - captions: Ground-truth captions; an integer array of shape (N, T) where
87       each element is in the range  $0 \leq y[i, t] < V$ 
88
89     Returns a tuple of:
90     - loss: Scalar loss
91     - grads: Dictionary of gradients parallel to self.params
92     """
93     # Cut captions into two pieces: captions_in has everything but the last word
94     # and will be input to the RNN; captions_out has everything but the first
95     # word and this is what we will expect the RNN to generate. These are offset
96     # by one relative to each other because the RNN should produce word (t+1)
97     # after receiving word t. The first element of captions_in will be the START
98     # token, and the first element of captions_out will be the first word.
99     captions_in = captions[:, :-1]
100    captions_out = captions[:, 1:]
101
102    # You'll need this
103    mask = (captions_out != self._null)
104
105    # Weight and bias for the affine transform from image features to initial
106    # hidden state
107    W_proj, b_proj = self.params['W_proj'], self.params['b_proj']
108
109    # Word embedding matrix
110    W_embed = self.params['W_embed']
111
112    # Input-to-hidden, hidden-to-hidden, and biases for the RNN
113    Wx, Wh, b = self.params['Wx'], self.params['Wh'], self.params['b']
114
115    # Weight and bias for the hidden-to-vocab transformation.
116    W_vocab, b_vocab = self.params['W_vocab'], self.params['b_vocab']
117
118    loss, grads = 0.0, {}
119    #####
120    # TODO: Implement the forward and backward passes for the CaptioningRNN.
121    # In the forward pass you will need to do the following:
122    # (1) Use an affine transformation to compute the initial hidden state
123    # from the image features. This should produce an array of shape (N, H)
124    # (2) Use a word embedding layer to transform the words in captions_in
125    # from indices to vectors, giving an array of shape (N, T, W).
126    # (3) Use either a vanilla RNN or LSTM (depending on self.cell_type) to
127    # process the sequence of input word vectors and produce hidden state
128    # vectors for all timesteps, producing an array of shape (N, T, H).
129    # (4) Use a (temporal) affine transformation to compute scores over the
130    # vocabulary at every timestep using the hidden states, giving an
131    # array of shape (N, T, V).
132    # (5) Use (temporal) softmax to compute loss using captions_out, ignoring
133    # the points where the output word is <NULL> using the mask above.
134    #
135    # In the backward pass you will need to compute the gradient of the loss
136    # with respect to all model parameters. Use the loss and grads variables
137    # defined above to store loss and gradients; grads[k] should give the
138    # gradients for self.params[k].
139    #
140    # Note also that you are allowed to make use of functions from layers.py
141    # in your implementation, if needed.
142    #####
143    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
144
145    # Forward pass
146    initial_hidden, initial_cache = affine_forward(features, W_proj, b_proj)
147    captions_embedded, embedding_cache = word_embedding_forward(captions_in, W_embed)
148    if self.cell_type == "rnn":
149        rnn_output, rnn_cache = rnn_forward(captions_embedded, initial_hidden, Wx, Wh, b)

```

```

150     elif self.cell_type == 'lstm':
151         rnn_output, lstm_cache = lstm_forward(captions_embedded, initial_hidden, Wx, Wh, b)
152
153     scores, vocab_cache = temporal_affine_forward(rnn_output, W_vocab, b_vocab)
154     loss, dx = temporal_softmax_loss(scores, captions_out, mask)
155
156     # Backward pass
157     dscores, grads['W_vocab'], grads['b_vocab'] = temporal_affine_backward(dx, vocab_cache)
158
159     if self.cell_type == 'rnn':
160         dx, dh_initial, grads['Wx'], grads['Wh'], grads['b'] = rnn_backward(dscores, rnn_cache)
161     elif self.cell_type == 'lstm':
162         dx, dh_initial, grads['Wx'], grads['Wh'], grads['b'] = lstm_backward(dscores, lstm_cache)
163
164     grads['W_embed'] = word_embedding_backward(dx, embedding_cache)
165     dx_initial, grads['W_proj'], grads['b_proj'] = affine_backward(dh_initial, initial_cache)
166
167     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
168     #####
169     #                                     END OF YOUR CODE                                     #
170     #####
171
172     return loss, grads
173
174
175 def sample(self, features, max_length=30):
176     """
177     Run a test-time forward pass for the model, sampling captions for input
178     feature vectors.
179
180     At each timestep, we embed the current word, pass it and the previous hidden
181     state to the RNN to get the next hidden state, use the hidden state to get
182     scores for all vocab words, and choose the word with the highest score as
183     the next word. The initial hidden state is computed by applying an affine
184     transform to the input image features, and the initial word is the <START>
185     token.
186
187     For LSTMs you will also have to keep track of the cell state; in that case
188     the initial cell state should be zero.
189
190     Inputs:
191     - features: Array of input image features of shape (N, D).
192     - max_length: Maximum length T of generated captions.
193
194     Returns:
195     - captions: Array of shape (N, max_length) giving sampled captions,
196       where each element is an integer in the range [0, V). The first element
197       of captions should be the first sampled word, not the <START> token.
198     """
199     N = features.shape[0]
200     captions = self._null * np.ones((N, max_length), dtype=np.int32)
201
202     # Unpack parameters
203     W_proj, b_proj = self.params['W_proj'], self.params['b_proj']
204     W_embed = self.params['W_embed']
205     Wx, Wh, b = self.params['Wx'], self.params['Wh'], self.params['b']
206     W_vocab, b_vocab = self.params['W_vocab'], self.params['b_vocab']
207
208     #####
209     # TODO: Implement test-time sampling for the model. You will need to
210     # initialize the hidden state of the RNN by applying the learned affine
211     # transform to the input image features. The first word that you feed to
212     # the RNN should be the <START> token; its value is stored in the
213     # variable self._start. At each timestep you will need to do to:
214     # (1) Embed the previous word using the learned word embeddings
215     # (2) Make an RNN step using the previous hidden state and the embedded
216     #     current word to get the next hidden state.
217     # (3) Apply the learned affine transformation to the next hidden state to
218     #     get scores for all words in the vocabulary
219     # (4) Select the word with the highest score as the next word, writing it
220     #     (the word index) to the appropriate slot in the captions variable
221     #
222     # For simplicity, you do not need to stop generating after an <END> token
223     # is sampled, but you can if you want to.
224     #
225     # HINT: You will not be able to use the rnn_forward or lstm_forward

```

```

226 # functions; you'll need to call rnn_step_forward or lstm_step_forward in #
227 # a loop.                                                                    #
228 #                                                                            #
229 # NOTE: we are still working over minibatches in this function. Also if    #
230 # you are using an LSTM, initialize the first cell state to zeros.         #
231 #####                                                                    #
232 # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****        #
233
234 hidden_state, hidden_cache = affine_forward(features, W_proj, b_proj)
235
236 if self.cell_type == 'lstm':
237     cell_state = np.zeros_like(hidden_state)
238
239 word_embed, embedding_cache = word_embedding_forward(self._start, W_embed)
240
241 for i in range(max_length):
242     if self.cell_type == 'rnn':
243         hidden_state, rnn_cache = rnn_step_forward(word_embed, hidden_state, Wx, Wh, b)
244     elif self.cell_type == 'lstm':
245         hidden_state, cell_state, rnn_cache = lstm_step_forward(word_embed, hidden_state,
cell_state, Wx, Wh, b)
246
247         scores, affine_cache = affine_forward(hidden_state, W_vocab, b_vocab)
248
249         captions[:, i] = np.argmax(scores, axis=1)
250
251         word_embed, _ = word_embedding_forward(captions[:, i], W_embed)
252
253 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
254 #####                                                                    #
255 #                               END OF YOUR CODE                               #
256 #####                                                                    #
257 return captions

```

2 rnn_layers.py

```
1 from future import print_function, division
2 from builtins import range
3 import numpy as np
4
5
6 """
7 This file defines layer types that are commonly used for recurrent neural
8 networks.
9 """
10
11
12 def rnn_step_forward(x, prev_h, Wx, Wh, b):
13     """
14     Run the forward pass for a single timestep of a vanilla RNN that uses a tanh
15     activation function.
16
17     The input data has dimension D, the hidden state has dimension H, and we use
18     a minibatch size of N.
19
20     Inputs:
21     - x: Input data for this timestep, of shape (N, D).
22     - prev_h: Hidden state from previous timestep, of shape (N, H)
23     - Wx: Weight matrix for input-to-hidden connections, of shape (D, H)
24     - Wh: Weight matrix for hidden-to-hidden connections, of shape (H, H)
25     - b: Biases of shape (H,)
26
27     Returns a tuple of:
28     - next_h: Next hidden state, of shape (N, H)
29     - cache: Tuple of values needed for the backward pass.
30     """
31     next_h, cache = None, None
32     #####
33     # TODO: Implement a single forward step for the vanilla RNN. Store the next #
34     # hidden state and any values you need for the backward pass in the next_h #
35     # and cache variables respectively. #
36     #####
37     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
38
39     next_h = np.tanh(np.dot(x, Wx) + np.dot(prev_h, Wh) + b)
40     cache = (x, prev_h, Wx, Wh, b, next_h)
41
42     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
43     #####
44     #                               END OF YOUR CODE                               #
45     #####
46     return next_h, cache
47
48
49 def rnn_step_backward(dnext_h, cache):
50     """
51     Backward pass for a single timestep of a vanilla RNN.
52
53     Inputs:
54     - dnext_h: Gradient of loss with respect to next hidden state, of shape (N, H)
55     - cache: Cache object from the forward pass
56
57     Returns a tuple of:
58     - dx: Gradients of input data, of shape (N, D)
59     - dprev_h: Gradients of previous hidden state, of shape (N, H)
60     - dWx: Gradients of input-to-hidden weights, of shape (D, H)
61     - dWh: Gradients of hidden-to-hidden weights, of shape (H, H)
62     - db: Gradients of bias vector, of shape (H,)
63     """
64     dx, dprev_h, dWx, dWh, db = None, None, None, None, None
65     #####
66     # TODO: Implement the backward pass for a single step of a vanilla RNN. #
67     # # #
68     # HINT: For the tanh function, you can compute the local derivative in terms #
69     # of the output value from tanh. #
70     #####
71     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
72
73     x, prev_h, Wx, Wh, b, next_h = cache
```

```

74 dhraw = (1 - next_h * next_h) * dnext_h
75 db = np.sum(dhraw, axis = 0)
76 dx = np.dot(dhraw, Wx.T)
77 dWx = np.dot(x.T, dhraw)
78 dprev_h = np.dot(dhraw, Wh.T)
79 dWh = np.dot(prev_h.T, dhraw)
80
81
82
83 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
84 #####
85 #                                     END OF YOUR CODE                                     #
86 #####
87 return dx, dprev_h, dWx, dWh, db
88
89
90 def rnn_forward(x, h0, Wx, Wh, b):
91     """
92     Run a vanilla RNN forward on an entire sequence of data. We assume an input
93     sequence composed of T vectors, each of dimension D. The RNN uses a hidden
94     size of H, and we work over a minibatch containing N sequences. After running
95     the RNN forward, we return the hidden states for all timesteps.
96
97     Inputs:
98     - x: Input data for the entire timeseries, of shape (N, T, D).
99     - h0: Initial hidden state, of shape (N, H)
100    - Wx: Weight matrix for input-to-hidden connections, of shape (D, H)
101    - Wh: Weight matrix for hidden-to-hidden connections, of shape (H, H)
102    - b: Biases of shape (H,)
103
104    Returns a tuple of:
105    - h: Hidden states for the entire timeseries, of shape (N, T, H).
106    - cache: Values needed in the backward pass
107    """
108    h, cache = None, None
109    #####
110    # TODO: Implement forward pass for a vanilla RNN running on a sequence of #
111    # input data. You should use the rnn_step_forward function that you defined #
112    # above. You can use a for loop to help compute the forward pass.         #
113    #####
114    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
115
116    N, T, D = x.shape
117    H = h0.shape[1]
118
119    h = np.zeros([N, T, H])
120    cache = []
121
122
123    for i in range(T):
124        next_h, cache_i = rnn_step_forward(x[:, i, :], h0, Wx, Wh, b)
125        h0 = next_h
126        h[:, i, :] = next_h
127        cache.append(cache_i)
128
129    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
130    #####
131    #                                     END OF YOUR CODE                                     #
132    #####
133    return h, cache
134
135
136 def rnn_backward(dh, cache):
137     """
138     Compute the backward pass for a vanilla RNN over an entire sequence of data.
139
140     Inputs:
141     - dh: Upstream gradients of all hidden states, of shape (N, T, H).
142
143     NOTE: 'dh' contains the upstream gradients produced by the
144     individual loss functions at each timestep, *not* the gradients
145     being passed between timesteps (which you'll have to compute yourself
146     by calling rnn_step_backward in a loop).
147
148     Returns a tuple of:
149     - dx: Gradient of inputs, of shape (N, T, D)

```



```

226
227 def word_embedding_backward(dout, cache):
228     """
229     Backward pass for word embeddings. We cannot back-propagate into the words
230     since they are integers, so we only return gradient for the word embedding
231     matrix.
232
233     HINT: Look up the function np.add.at
234
235     Inputs:
236     - dout: Upstream gradients of shape (N, T, D)
237     - cache: Values from the forward pass
238
239     Returns:
240     - dW: Gradient of word embedding matrix, of shape (V, D).
241     """
242     dW = None
243     #####
244     # TODO: Implement the backward pass for word embeddings.
245     #
246     # Note that words can appear more than once in a sequence.
247     # HINT: Look up the function np.add.at
248     #####
249     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
250
251     x, W = cache
252     V = W.shape[0]
253     dW = np.zeros_like(W)
254
255     # Add dout of the corresponding indices (x) to dW
256     np.add.at(dW, x, dout)
257
258     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
259     #####
260     #                               END OF YOUR CODE
261     #####
262     return dW
263
264
265 def sigmoid(x):
266     """
267     A numerically stable version of the logistic sigmoid function.
268     """
269     pos_mask = (x >= 0)
270     neg_mask = (x < 0)
271     z = np.zeros_like(x)
272     z[pos_mask] = np.exp(-x[pos_mask])
273     z[neg_mask] = np.exp(x[neg_mask])
274     top = np.ones_like(x)
275     top[neg_mask] = z[neg_mask]
276     return top / (1 + z)
277
278
279 def lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b):
280     """
281     Forward pass for a single timestep of an LSTM.
282
283     The input data has dimension D, the hidden state has dimension H, and we use
284     a minibatch size of N.
285
286     Note that a sigmoid() function has already been provided for you in this file.
287
288     Inputs:
289     - x: Input data, of shape (N, D)
290     - prev_h: Previous hidden state, of shape (N, H)
291     - prev_c: previous cell state, of shape (N, H)
292     - Wx: Input-to-hidden weights, of shape (D, 4H)
293     - Wh: Hidden-to-hidden weights, of shape (H, 4H)
294     - b: Biases, of shape (4H,)
295
296     Returns a tuple of:
297     - next_h: Next hidden state, of shape (N, H)
298     - next_c: Next cell state, of shape (N, H)
299     - cache: Tuple of values needed for backward pass.
300     """
301     next_h, next_c, cache = None, None, None

```



```

378 dWx = np.dot(x.T, dA)
379 dprev_h = np.dot(dA, Wh.T)
380 dWh = np.dot(prev_h.T, dA)
381 db = np.sum(dA, axis=0)
382
383
384
385 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
386 #####
387 #                                     END OF YOUR CODE                                     #
388 #####
389
390 return dx, dprev_h, dprev_c, dWx, dWh, db
391
392
393 def lstm_forward(x, h0, Wx, Wh, b):
394     """
395     Forward pass for an LSTM over an entire sequence of data. We assume an input
396     sequence composed of T vectors, each of dimension D. The LSTM uses a hidden
397     size of H, and we work over a minibatch containing N sequences. After running
398     the LSTM forward, we return the hidden states for all timesteps.
399
400     Note that the initial cell state is passed as input, but the initial cell
401     state is set to zero. Also note that the cell state is not returned; it is
402     an internal variable to the LSTM and is not accessed from outside.
403
404     Inputs:
405     - x: Input data of shape (N, T, D)
406     - h0: Initial hidden state of shape (N, H)
407     - Wx: Weights for input-to-hidden connections, of shape (D, 4H)
408     - Wh: Weights for hidden-to-hidden connections, of shape (H, 4H)
409     - b: Biases of shape (4H,)
410
411     Returns a tuple of:
412     - h: Hidden states for all timesteps of all sequences, of shape (N, T, H)
413     - cache: Values needed for the backward pass.
414     """
415     h, cache = None, None
416     #####
417     # TODO: Implement the forward pass for an LSTM over an entire timeseries. #
418     # You should use the lstm_step_forward function that you just defined. #
419     #####
420     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
421
422     N, T, D = x.shape
423     H = h0.shape[1]
424
425     prev_h = h0
426     prev_c = np.zeros_like(h0)
427
428     h = np.zeros([N, T, H])
429     cache = []
430
431
432     for i in range(T):
433         next_h, next_c, cache_i = lstm_step_forward(x[:, i, :], prev_h, prev_c, Wx, Wh, b)
434         prev_h = next_h
435         prev_c = next_c
436         h[:, i, :] = next_h
437         cache.append(cache_i)
438
439     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
440     #####
441     #                                     END OF YOUR CODE                                     #
442     #####
443
444     return h, cache
445
446
447 def lstm_backward(dh, cache):
448     """
449     Backward pass for an LSTM over an entire sequence of data.]
450
451     Inputs:
452     - dh: Upstream gradients of hidden states, of shape (N, T, H)
453     - cache: Values from the forward pass

```

```

454 Returns a tuple of:
455 - dx: Gradient of input data of shape (N, T, D)
456 - dh0: Gradient of initial hidden state of shape (N, H)
457 - dWx: Gradient of input-to-hidden weight matrix of shape (D, 4H)
458 - dWh: Gradient of hidden-to-hidden weight matrix of shape (H, 4H)
459 - db: Gradient of biases, of shape (4H,)
460 """
461 dx, dh0, dWx, dWh, db = None, None, None, None, None
462 #####
463 # TODO: Implement the backward pass for an LSTM over an entire timeseries. #
464 # You should use the lstm_step_backward function that you just defined. #
465 #####
466 # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
467
468 x, prev_h, prev_c, Wx, Wh, b, i, f, o, g, next_c, next_h = cache[0]
469 N, T, H = dh.shape
470 D = Wx.shape[0]
471
472 dx = np.zeros([N, T, D])
473 dWx = np.zeros_like(Wx)
474 dWh = np.zeros_like(Wh)
475 db = np.zeros_like(b)
476 dprev_h = np.zeros_like(prev_h)
477 dprev_c = np.zeros_like(prev_c)
478
479 dh0 = np.zeros([N, H])
480
481 for i in reversed(range(T)):
482     dh_i = dprev_h + dh[:, i, :]
483
484     dx[:, i, :], dprev_h, dprev_c, dWx_i, dWh_i, db_i = lstm_step_backward(dh_i, dprev_c, cache[i])
485
486     db += db_i
487     dWh += dWh_i
488     dWx += dWx_i
489
490
491 dh0 = dprev_h
492
493
494
495 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
496 #####
497 # END OF YOUR CODE #
498 #####
499
500 return dx, dh0, dWx, dWh, db
501
502
503
504 def temporal_affine_forward(x, w, b):
505     """
506     Forward pass for a temporal affine layer. The input is a set of D-dimensional
507     vectors arranged into a minibatch of N timeseries, each of length T. We use
508     an affine function to transform each of those vectors into a new vector of
509     dimension M.
510
511     Inputs:
512     - x: Input data of shape (N, T, D)
513     - w: Weights of shape (D, M)
514     - b: Biases of shape (M,)
515
516     Returns a tuple of:
517     - out: Output data of shape (N, T, M)
518     - cache: Values needed for the backward pass
519     """
520     N, T, D = x.shape
521     M = b.shape[0]
522     out = x.reshape(N * T, D).dot(w).reshape(N, T, M) + b
523     cache = x, w, b, out
524     return out, cache
525
526
527 def temporal_affine_backward(dout, cache):
528     """
529     Backward pass for temporal affine layer.

```

```

530
531     Input:
532     - dout: Upstream gradients of shape (N, T, M)
533     - cache: Values from forward pass
534
535     Returns a tuple of:
536     - dx: Gradient of input, of shape (N, T, D)
537     - dw: Gradient of weights, of shape (D, M)
538     - db: Gradient of biases, of shape (M,)
539     """
540     x, w, b, out = cache
541     N, T, D = x.shape
542     M = b.shape[0]
543
544     dx = dout.reshape(N * T, M).dot(w.T).reshape(N, T, D)
545     dw = dout.reshape(N * T, M).T.dot(x.reshape(N * T, D)).T
546     db = dout.sum(axis=(0, 1))
547
548     return dx, dw, db
549
550
551 def temporal_softmax_loss(x, y, mask, verbose=False):
552     """
553     A temporal version of softmax loss for use in RNNs. We assume that we are
554     making predictions over a vocabulary of size V for each timestep of a
555     timeseries of length T, over a minibatch of size N. The input x gives scores
556     for all vocabulary elements at all timesteps, and y gives the indices of the
557     ground-truth element at each timestep. We use a cross-entropy loss at each
558     timestep, summing the loss over all timesteps and averaging across the
559     minibatch.
560
561     As an additional complication, we may want to ignore the model output at some
562     timesteps, since sequences of different length may have been combined into a
563     minibatch and padded with NULL tokens. The optional mask argument tells us
564     which elements should contribute to the loss.
565
566     Inputs:
567     - x: Input scores, of shape (N, T, V)
568     - y: Ground-truth indices, of shape (N, T) where each element is in the range
569         0 <= y[i, t] < V
570     - mask: Boolean array of shape (N, T) where mask[i, t] tells whether or not
571         the scores at x[i, t] should contribute to the loss.
572
573     Returns a tuple of:
574     - loss: Scalar giving loss
575     - dx: Gradient of loss with respect to scores x.
576     """
577
578     N, T, V = x.shape
579
580     x_flat = x.reshape(N * T, V)
581     y_flat = y.reshape(N * T)
582     mask_flat = mask.reshape(N * T)
583
584     probs = np.exp(x_flat - np.max(x_flat, axis=1, keepdims=True))
585     probs /= np.sum(probs, axis=1, keepdims=True)
586     loss = -np.sum(mask_flat * np.log(probs[np.arange(N * T), y_flat])) / N
587     dx_flat = probs.copy()
588     dx_flat[np.arange(N * T), y_flat] -= 1
589     dx_flat /= N
590     dx_flat *= mask_flat[:, None]
591
592     if verbose: print('dx_flat: ', dx_flat.shape)
593
594     dx = dx_flat.reshape(N, T, V)
595
596     return loss, dx

```