
A Tutorial On Backward Propagation Through Time (BPTT) In The Gated Recurrent Unit (GRU) RNN

Minchen Li

Department of Computer Science
The University of British Columbia
minchenl@cs.ubc.ca

Abstract

In this tutorial, we provide a thorough explanation on how BPTT in GRU¹ is conducted. A MATLAB program which implements the entire BPTT for GRU and the pseudo-codes describing the algorithms explicitly will be presented. We provide two algorithms for BPTT, a direct but quadratic time algorithm for easy understanding, and an optimized linear time algorithm. This tutorial starts with a specification of the problem followed by a mathematical derivation before the computational solutions.

1 Specification

We want to use a dataset containing n_s sentences each with n_w words to train a GRU language model, and our vocabulary size is n_v . Namely, we have input $x \in R^{n_v \times n_w \times n_s}$ and label $y \in R^{n_v \times n_w \times n_s}$ both representing n_s sentences.

For simplicity, let's look at one sentence at a time. In one sentence, the one-hot vector $x_t \in R^{n_v \times 1}$ represents the t^{th} word. For time step t , the GRU unit computes the output \hat{y}_t using the input x_t and the previous internal state s_{t-1} as follows:

$$\begin{aligned}z_t &= \sigma(U_z x_t + W_z s_{t-1} + b_z) \\r_t &= \sigma(U_r x_t + W_r s_{t-1} + b_r) \\h_t &= \tanh(U_h x_t + W_h (s_{t-1} \odot r_t) + b_h) \\s_t &= (1 - z_t) \odot h_t + z_t \odot s_{t-1} \\ \hat{y}_t &= \text{softmax}(V s_t + b_V)\end{aligned}\tag{1}$$

Here \odot is the vector element-wise multiplication, $\sigma(\cdot)$ is the element-wise *sigmoid* function, and $\tanh(\cdot)$ is the element-wise *hyperbolic tangent* function. The dimensions of the parameters are as follows:

$$\begin{aligned}U_z, U_r, U_h &\in R^{n_i \times n_v} \\W_z, W_r, W_h &\in R^{n_i \times n_i} \\b_z, b_r, b_h &\in R^{n_i \times 1} \\V &\in R^{n_v \times n_i}, b_V \in R^{n_v \times 1}\end{aligned}$$

where n_i is the internal memory size set by the user.

¹GRU is an improved version of traditional RNN (Recurrent Neural Network, see WildML.com for an introduction. This link also provides an introduction to GRU and some general discussion on BPTT and beyond.)

Then for step t , we can calculate the cross entropy loss L_t as:

$$L_t = \text{sumOfAllElements} \left(-y_t \odot \log(\hat{y}_t) \right) \quad (2)$$

Here \log is also an element-wise function.

To train the GRU, we want to know the values of all parameters that minimize the total loss $L = \sum_{t=1}^{n_w} L_t$:

$$\underset{\Theta}{\operatorname{argmin}} L$$

where $\Theta = \{U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V\}$. This is a non-convex problem with huge input data. So people usually use Stochastic Gradient Descent² method to solve this problem, which means we need to calculate $\partial L / \partial U_z, \partial L / \partial U_r, \partial L / \partial U_h, \partial L / \partial W_z, \partial L / \partial W_r, \partial L / \partial W_h, \partial L / \partial b_z, \partial L / \partial b_r, \partial L / \partial b_h, \partial L / \partial V, \partial L / \partial b_V$ given a batch of sentences. (Note that in each step, these parameters stays the same.) In this tutorial we consider using only one sentence at a time to make it concise.

2 Derivation

The best way to calculate gradients using the Chain Rule from output to input is to first draw the expression graph of the entire model in order to figure out the relations between the output, intermediate results, and the input³. Here we draw part of the expression graph of GRU in Fig.1.

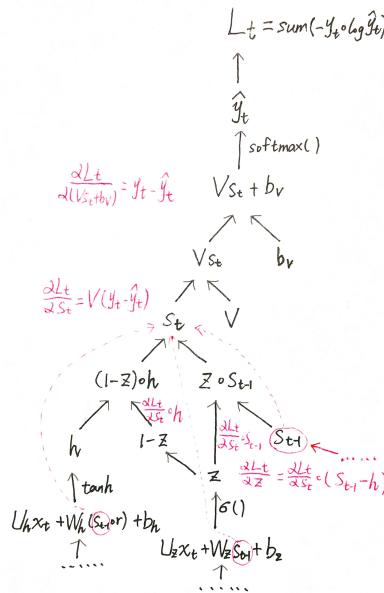


Figure 1: The upper part of expression graph describing the operations of GRU. Note that the sub-graph which s_{t-1} depends on is just like the sub-graph of s_t . This is what the red dashed lines mean.

With this expression graph, the Chain Rule works if you go backwards along the edges (top-down). If a node X has multiple outgoing edges connecting the target node T , you need to sum over the partial derivatives of each of those outgoing edges to derive the gradient $\partial T / \partial X$. We will illustrate the rules in the following paragraphs.

Let's take $\partial L / \partial U_z$ as the example here. Others are just similar. Since $L = \sum_{t=1}^{n_w} L_t$ and the parameters stay the same in each step, we also have $\partial L / \partial U_z = \sum_{t=1}^{n_w} (\partial L_t / \partial U_z)$, so let's calculate each $\partial L_t / \partial U_z$ independently and sum them up.

²See the Wikipedia to get some knowledge about Stochastic Gradient Descent.

³See colah's blog and Stanford CS231n Course Note for some general introductions.

With the Chain Rule, we have:

$$\frac{\partial L_t}{\partial U_z} = \frac{\partial L_t}{\partial s_t} \frac{\partial s_t}{\partial U_z} \quad (3)$$

The first part is just trivial if you know how to differentiate the cross entropy loss function embedded with the *softmax* function:

$$\frac{\partial L_t}{\partial s_t} = V(\hat{y}_t - y_t)$$

For $\partial z/\partial U_z$, similarly, some people might just derive: (if they know how to differentiate *sigmoid* function)

$$\overline{\frac{\partial s_t}{\partial U_z}} = \left((s_{t-1} - h_t) \odot z_t \odot (1 - z_t) \right) x_t^T \quad (4)$$

Here there are two expressions $1 - z$ and $z \odot s_{t-1}$ influencing $\partial s_t/\partial z$ as shown in our expression graph. The solution is to derive partial derivatives through each edge and then add them up, which is exactly how we deal with $\partial s_t/\partial s_{t-1}$ as you will see in the following paragraphs. However, Eq.4 only calculates one part of the gradient, so we put a bar on top of it, while you may find this very useful in our following calculations.

Note that s_{t-1} also depends on U_z , so we can not treat it as a constant here. Moreover, this s_{t-1} will also introduce the influence of s_i , where $i = 1, \dots, t-2$. So for clearness, we should expand Eq.3 as:

$$\begin{aligned} \frac{\partial L_t}{\partial U_z} &= \frac{\partial L_t}{\partial s_t} \frac{\partial s_t}{\partial U_z} \\ &= \frac{\partial L_t}{\partial s_t} \sum_{i=1}^t \left(\frac{\partial s_t}{\partial s_i} \overline{\frac{\partial s_i}{\partial U_z}} \right) \\ &= \frac{\partial L_t}{\partial s_t} \sum_{i=1}^t \left(\left(\prod_{j=i}^{t-1} \frac{\partial s_{j+1}}{\partial s_j} \right) \overline{\frac{\partial s_i}{\partial U_z}} \right) \end{aligned} \quad (5)$$

where $\overline{\partial s_i/\partial U_z}$ is the gradient of s_i with respect to U_z while taking s_{i-1} as a constant, of which a similar example has been shown in Eq.4 for step t .

The derivation of $\partial s_t/\partial s_{t-1}$ is similar to the derivation of $\partial s_t/\partial z$ as has been discussed above. Since there are four outgoing edges from s_{t-1} to s_t directly and indirectly through z_t , r_t , and h_t in the expression graph, we need to sum all the four partial derivatives together:

$$\begin{aligned} \frac{\partial s_t}{\partial s_{t-1}} &= \frac{\partial s_t}{\partial h_t} \frac{\partial h_t}{\partial s_{t-1}} + \frac{\partial s_t}{\partial z_t} \frac{\partial z_t}{\partial s_{t-1}} + \overline{\frac{\partial s_t}{\partial s_{t-1}}} \\ &= \frac{\partial s_t}{\partial h_t} \left(\frac{\partial h_t}{\partial r_t} \frac{\partial r_t}{\partial s_{t-1}} + \overline{\frac{\partial h_t}{\partial s_{t-1}}} \right) + \frac{\partial s_t}{\partial z_t} \frac{\partial z_t}{\partial s_{t-1}} + \overline{\frac{\partial s_t}{\partial s_{t-1}}} \end{aligned} \quad (6)$$

where $\overline{\partial s_t/\partial s_{t-1}}$ is the gradient of s_t with respect to s_{t-1} while taking h_t and z_t as constants. Similarly, $\overline{\partial h_t/\partial s_{t-1}}$ is the gradient of h_t with respect to s_{t-1} while taking r_t as a constant.

Plugging the intermediate results in the above formula, we get:

$$\begin{aligned} \frac{\partial s_t}{\partial s_{t-1}} &= (1 - z_t) \left(W_r^T \left((W_h^T (1 - h \odot h)) \odot s_{t-1} \odot r \odot (1 - r) \right) + \left((W_h^T (1 - h \odot h)) \odot r_t \right) + \right. \\ &\quad \left. W_z^T \left((s_{t-1} - h_t) \odot z_t \odot (1 - z_t) \right) \right) + z \end{aligned}$$

Till now, we have covered all the components needed to calculate $\partial L_t/\partial U_z$. The gradient of L_t with respect to other parameters are just similar. In the next chapter, we will provide a more machinery view of the calculation - the pseudo-code describing the algorithm to calculate the gradients. In the last chapter of this tutorial, we will provide the pure machine representation - a MATLAB program which implements the calculation and verification of BPTT. If you just want to understand the idea behind BPTT and decide to use fully supported auto-differentiation packages (like Theano⁴) to build your own GRU, you can stop here. If you need to implement the exact chain rule like us or just curious about what will happen next, get ready to proceed!

⁴ Theano is a Python library that allows you to define, optimize, and evaluate mathematical expressions involving multi-dimensional arrays efficiently.

3 Algorithm

Here we also only take $\partial L/\partial U_z$ as the example. We will provide the calculation of all the gradients in the next chapter.

We present two algorithms, one direct algorithm as derived previously calculating $\partial L_t/\partial U_z$ and sum them up while taking $O(n_w^2)$ time, and the other $O(n_w)$ time algorithm which we will see later.

Algorithm 1 A direct but $O(n_w^2)$ time algorithm to calculate $\partial L/\partial U_z$ (and beyond)

Input: The training data $X, Y \in R^{n_v \times n_w}$ composed of the one-hot column vectors $x_t, y_t \in R^{n_v \times 1}$, $t = 1, 2, \dots, n_w$ representing the words in the sentence.

Input: A vector $s_0 \in R^{n_i \times 1}$ representing the initial internal state of the model (usually set to 0).

Input: The parameters $\Theta = \{U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V\}$ of the model.

Output: The total loss gradient $\partial L/\partial U_z$.

- 1: %forward propagate to calculate the internal states $S \in R^{n_i \times n_w}$, the predictions $\hat{Y} \in R^{n_v \times n_w}$, the losses $L_{mtr} \in R^{n_w \times 1}$, and the intermediate results $Z, R, C \in R^{n_i \times n_w}$ of each step:
 - 2: $[S, \hat{Y}, L_{mtr}, Z, R, C] = \text{forward}(X, Y, \Theta, s_0)$ % forward() can be implemented easily according to Eq.1 and Eq.2
 - 3: $dU_z = \text{zeros}(n_i, n_w)$ % initialize a variable dU_z
 - 4: $\partial L_{mtr}/\partial S = V^T(\hat{Y} - Y)$ % calculate $\partial L_t/\partial s_t$ for $t = 1, 2, \dots, n_w$ with one matrix operation
 - 5: **for** $t \leftarrow 1$ **to** n_w % calculate each $\partial L_t/\partial U_z$ and accumulate
 - 6: **for** $j \leftarrow t$ **to** 1 % calculate each $(\partial L_t/\partial s_j)(\overline{\partial s_j/\partial U_z})$ and accumulate
 - 7: $\partial L_t/\partial z_j = \partial L_t/\partial s_j \odot (s_{j-1} - h_j)$ % $\partial s_j/\partial z_j$ is $(s_{j-1} - h_j)$, $\partial L_t/\partial s_j$ is calculated in the last inner loop iteration or in Line 4
 - 8: $\partial L_t/\partial(U_z x_j + W_z s_{j-1} + b_z) = \partial L_t/\partial z_j \odot z_j \odot (1 - z_j)$ % $\partial \sigma(x)/\partial x = \sigma(x) \odot (1 - \sigma(x))$
 - 9: $dU_z += (\partial L_t/\partial(U_z x_j + W_z s_{j-1} + b_z))x_j^T$ % accumulate
 - 10: calculate $\partial L_t/\partial s_{j-1}$ using $\partial L_t/\partial s_j$ and Eq.6 % for the next inner loop iteration
 - 11: **end**
 - 12: **end**
 - 13: **return** dU_z % $\partial L/\partial U_z$
-

The above direct algorithm actually follows Eq.5 to calculate $\partial L_t/\partial U_z$ and then add them up to form $\partial L/\partial U_z$:

$$\begin{aligned} \frac{\partial L}{\partial U_z} &= \sum_{t=1}^{n_w} \frac{\partial L_t}{\partial U_z} \\ &= \sum_{t=1}^{n_w} \left(\frac{\partial L_t}{\partial s_t} \sum_{i=1}^t \left(\frac{\partial s_t}{\partial s_i} \overline{\frac{\partial s_i}{\partial U_z}} \right) \right) \\ &= \sum_{t=1}^{n_w} \left(\frac{\partial L_t}{\partial s_t} \sum_{i=1}^t \left(\left(\prod_{j=i}^{t-1} \frac{\partial s_{j+1}}{\partial s_j} \right) \overline{\frac{\partial s_i}{\partial U_z}} \right) \right) \end{aligned}$$

If we just expand $\partial L_t/\partial U_z$ to the second line of the above equation and do some reordering, we can get:

$$\begin{aligned} \frac{\partial L}{\partial U_z} &= \sum_{t=1}^{n_w} \left(\frac{\partial L_t}{\partial s_t} \sum_{i=1}^t \left(\frac{\partial s_t}{\partial s_i} \overline{\frac{\partial s_i}{\partial U_z}} \right) \right) \\ &= \sum_{t=1}^{n_w} \left(\sum_{i=1}^t \left(\frac{\partial L_t}{\partial s_t} \frac{\partial s_t}{\partial s_i} \overline{\frac{\partial s_i}{\partial U_z}} \right) \right) \\ &= \sum_{t=1}^{n_w} \left(\sum_{i=1}^t \left(\frac{\partial L_t}{\partial s_i} \overline{\frac{\partial s_i}{\partial U_z}} \right) \right) \end{aligned}$$

Right now the inner summation keeps the subscript of ∂L_t and iterate over ∂s_i . If we further expand the inner summation and then sort them to iterate over ∂L_i , we get:

$$\frac{\partial L}{\partial U_z} = \sum_{t=1}^{n_w} \left(\left(\sum_{i=t}^{n_w} \frac{\partial L_i}{\partial s_t} \right) \overline{\frac{\partial s_t}{\partial U_z}} \right) \quad (7)$$

For the inner summation of Eq.7, we have:

$$\begin{aligned} \sum_{i=t}^{n_w} \left(\frac{\partial L_i}{\partial s_t} \right) &= \left(\sum_{i=t+1}^{n_w} \left(\frac{\partial L_i}{\partial s_{t+1}} \frac{\partial s_{t+1}}{\partial s_t} \right) \right) + \frac{\partial L_t}{\partial s_t} \\ &= \left(\sum_{i=t+1}^{n_w} \frac{\partial L_i}{\partial s_{t+1}} \right) \frac{\partial s_{t+1}}{\partial s_t} + \frac{\partial L_t}{\partial s_t} \end{aligned} \quad (8)$$

This just gives us an updating formula to calculate this inner summation for each step t incrementally rather than executing another *for* loop, thus making it possible for us to implement an $O(n_w)$ time algorithm!

Algorithm 2 An optimized $O(n_w)$ time algorithm to calculate $\partial L/\partial U_z$ (and beyond)

Input: The training data $X, Y \in R^{n_v \times n_w}$ composed of the one-hot column vectors $x_t, y_t \in R^{n_v \times 1}$, $t = 1, 2, \dots, n_w$ representing the words in the sentence.

Input: A vector $s_0 \in R^{n_i \times 1}$ representing the initial internal state of the model (usually set to 0).

Input: The parameters $\Theta = \{U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V\}$ of the model.

Output: The total loss gradient $\partial L/\partial U_z$.

- 1: %forward propagate to calculate the internal states $S \in R^{n_i \times n_w}$, the predictions $\hat{Y} \in R^{n_v \times n_w}$, the losses $L_{mtr} \in R^{n_w \times 1}$, and the intermediate results $Z, R, C \in R^{n_i \times n_w}$ of each step:
 - 2: $[S, \hat{Y}, L_{mtr}, Z, R, C] = \text{forward}(X, Y, \Theta, s_0)$ % forward() can be implemented easily according to Eq.1 and Eq.2
 - 3: $dU_z = \text{zeros}(n_i, n_w)$ % initialize a variable dU_z
 - 4: $\partial L_{mtr}/\partial S = V^T(\hat{Y} - Y)$ % calculate $\partial L_t/\partial s_t$ for $t = 1, 2, \dots, n_w$ with one matrix operation
 - 5: **for** $t \leftarrow n_w$ **to** 1 % calculate each $\left(\sum_{i=t}^{n_w} \left(\frac{\partial L_i}{\partial s_t} \right) \right) \overline{\frac{\partial s_t}{\partial U_z}}$ and accumulate
 - 6: $\sum_{i=t}^{n_w} (\partial L_i/\partial z_t) = \left(\sum_{i=t}^{n_w} (\partial L_i/\partial s_t) \right) \odot (s_{t-1} - h_t)$ % $\partial s_t/\partial z_t$ is $(s_{t-1} - h_t)$, $\sum_{i=t}^{n_w} (\partial L_i/\partial s_t)$ is calculated in the last iteration or in Line 4. (when $t = n_w$, $\sum_{i=t}^{n_w} (\partial L_i/\partial s_t) = \partial L_t/\partial s_t$)
 - 7: $\sum_{i=t}^{n_w} (\partial L_i/\partial (U_z x_t + W_z s_{t-1} + b_z)) = \left(\sum_{i=t}^{n_w} (\partial L_i/\partial z_t) \right) \odot z_t \odot (1 - z_t)$ % $\partial \sigma(x)/\partial x = \sigma(x) \odot (1 - \sigma(x))$
 - 8: $dU_z += \left(\sum_{i=t}^{n_w} (\partial L_i/\partial (U_z x_t + W_z s_{j-t} + b_z)) \right) x_t^T$ % accumulate
 - 9: calculate $\sum_{i=t-1}^{n_w} (\partial L_i/\partial s_{t-1})$ using Eq.6 and Eq.8 % for the next iteration
 - 10: **end**
 - 11: return dU_z % $\partial L/\partial U_z$
-

4 Implementation

Here we provide the MATLAB program which calculates the gradients with respect to all the parameters of GRU using our two proposed algorithms. It also checks the gradients with the numerical results. We will divide our code into two parts, the first part presented below contains the core functions implementing the BPTT of GRU we just derived, the second part is composed of some functions that are less important to the topic of this tutorial.

Core Functions

```
1 % This program tests the BPTT process we manually developed for GRU.
2 % We calculate the gradients of GRU parameters with chain rule, and then
3 % compare them to the numerical gradients to check whether our chain rule
4 % derivation is correct.
5
6 % Here, we provided 2 versions of BPTT, backward_direct() and backward().
7 % The former one is the direct idea to calculate gradient within each
8 % step
9 % and add them up (O(sentence_size^2) time). The latter one is optimized
10 % to
11 % calculate the contribution of each step to the overall gradient, which
12 % is
13 % only O(sentence_size) time.
14
15 % This is very helpful for people who wants to implement GRU in Caffe
16 % since
17 % Caffe didn't support auto-differentiation. This is also very helpful
18 % for
19 % the people who wants to know the details about Backpropagation Through
20 % Time algorithm in the Recurrent Neural Networks (such as GRU and LSTM)
21 % and also get a sense on how auto-differentiation is possible.
22
23 % NOTE: We didn't involve SGD training here. With SGD training, this
24 % program would become a complete implementation of GRU which can be
25 % trained with sequence data. However, since this is only a CPU serial
26 % Matlab version of GRU, applying it on large datasets will be
27 % dramatically
28 % slow.
29
30 % by Minchen Li, at The University of British Columbia. 2016-04-21
31
32 function testBPTT_GRU
33     % set GRU and data scale
34     vocabulary_size = 64;
35     iMem_size = 4;
36     sentence_size = 20; % number of words in a sentence
37                         % (including start and end symbol)
38                         % since we will only use one sentence for
39     training,          % this is also the total steps during training.
40
41     [x y] = getTrainingData(vocabulary_size, sentence_size);
42
43     % initialize parameters:
44     % multiplier for input x_t of intermediate variables
45     U_z = rand(iMem_size, vocabulary_size);
46     U_r = rand(iMem_size, vocabulary_size);
47     U_c = rand(iMem_size, vocabulary_size);
48     % multiplier for pervious s of intermediate variables
49     W_z = rand(iMem_size, iMem_size);
50     W_r = rand(iMem_size, iMem_size);
51     W_c = rand(iMem_size, iMem_size);
52     % bias terms of intermediate variables
53     b_z = rand(iMem_size, 1);
```

```

49     b_r = rand(iMem_size, 1);
    b_c = rand(iMem_size, 1);
    % decoder for generating output
51     V = rand(vocabulary_size, iMem_size);
    b_V = rand(vocabulary_size, 1); % bias of decoder
53     % previous s of step 1
    s_0 = rand(iMem_size, 1);
55
    % calculate and check gradient
57     tic
    [dV, db_V, dU_z, dU_r, dU_c, dW_z, dW_r, dW_c, db_z, db_r, db_c, ds_0
    ] = ...
59     backward_direct(x, y, U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c
    , V, b_V, s_0);
    toc
61     tic
    checkGradient_GRU(x, y, U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c,
    V, b_V, s_0, ...
63     dV, db_V, dU_z, dU_r, dU_c, dW_z, dW_r, dW_c, db_z, db_r, db_c,
    ds_0);
    toc
65
    tic
67     [dV, db_V, dU_z, dU_r, dU_c, dW_z, dW_r, dW_c, db_z, db_r, db_c, ds_0
    ] = ...
    backward(x, y, U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V,
    b_V, s_0);
69     toc
    tic
71     checkGradient_GRU(x, y, U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c,
    V, b_V, s_0, ...
    dV, db_V, dU_z, dU_r, dU_c, dW_z, dW_r, dW_c, db_z, db_r, db_c,
    ds_0);
73     toc
end
75
% Forward propagate calculate s, y_hat, loss and intermediate variables
for each step
77 function [s, y_hat, L, z, r, c] = forward(x, y, ...
    U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V, s_0)
79     % count sizes
    [vocabulary_size, sentence_size] = size(x);
81     iMem_size = size(V, 2);
83
    % initialize results
    s = zeros(iMem_size, sentence_size);
85     y_hat = zeros(vocabulary_size, sentence_size);
    L = zeros(sentence_size, 1);
87     z = zeros(iMem_size, sentence_size);
    r = zeros(iMem_size, sentence_size);
89     c = zeros(iMem_size, sentence_size);
91
    % calculate result for step 1 since s_0 is not in s
    z(:,1) = sigmoid(U_z*x(:,1) + W_z*s_0 + b_z);
93     r(:,1) = sigmoid(U_r*x(:,1) + W_r*s_0 + b_r);
    c(:,1) = tanh(U_c*x(:,1) + W_c*(s_0.*r(:,1)) + b_c);
95     s(:,1) = (1-z(:,1)).*c(:,1) + z(:,1).*s_0;
    y_hat(:,1) = softmax(V*s(:,1) + b_V);
97     L(1) = sum(-y(:,1).*log(y_hat(:,1)));
    % calculate results for step 2 - sentence_size similarly
99     for wordI = 2:sentence_size
        z(:,wordI) = sigmoid(U_z*x(:,wordI) + W_z*s(:,wordI-1) + b_z);
101        r(:,wordI) = sigmoid(U_r*x(:,wordI) + W_r*s(:,wordI-1) + b_r);
        c(:,wordI) = tanh(U_c*x(:,wordI) + W_c*(s(:,wordI-1).*r(:,wordI))
        + b_c);

```

```

103     s(:, wordI) = (1-z(:, wordI)).*c(:, wordI) + z(:, wordI).*s(:, wordI
-1);
104     y_hat(:, wordI) = softmax(V*s(:, wordI) + b_V);
105     L(wordI) = sum(-y(:, wordI).*log(y_hat(:, wordI)));
106     end
107 end
108
109 % Backward propagate to calculate gradient using chain rule
110 % (O(sentence_size) time)
111 function [dV, db_V, dU_z, dU_r, dU_c, dW_z, dW_r, dW_c, db_z, db_r, db_c,
ds_0] = ...
112     backward(x, y, U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V,
s_0)
113 % forward propagate to get the intermediate and output results
114 [s, y_hat, L, z, r, c] = forward(x, y, U_z, U_r, U_c, W_z, W_r, W_c,
...
115     b_z, b_r, b_c, V, b_V, s_0);
116 % count sentence size
117 [~, sentence_size] = size(x);
118
119 % calculate gradient using chain rule
120 delta_y = y_hat - y;
121 db_V = sum(delta_y, 2);
122
123 dV = zeros(size(V));
124 for wordI = 1:sentence_size
125     dV = dV + delta_y(:, wordI)*s(:, wordI)';
126 end
127
128 ds_0 = zeros(size(s_0));
129 dU_c = zeros(size(U_c));
130 dU_r = zeros(size(U_r));
131 dU_z = zeros(size(U_z));
132 dW_c = zeros(size(W_c));
133 dW_r = zeros(size(W_r));
134 dW_z = zeros(size(W_z));
135 db_z = zeros(size(b_z));
136 db_r = zeros(size(b_r));
137 db_c = zeros(size(b_c));
138 ds_single = V'*delta_y;
139 % calculate the derivative contribution of each step and add them up
140 ds_cur = zeros(size(ds_single,1), 1);
141 for wordJ = sentence_size:-1:2
142     ds_cur = ds_cur + ds_single(:, wordJ);
143     ds_cur_bk = ds_cur;
144
145     dtanhInput = (ds_cur.*(1-z(:, wordJ)).*(1-c(:, wordJ).*c(:, wordJ)))
;
146     db_c = db_c + dtanhInput;
147     dU_c = dU_c + dtanhInput*x(:, wordJ)'; %could be accelerated by
avoiding add 0
148     dW_c = dW_c + dtanhInput*(s(:, wordJ-1).*r(:, wordJ))';
149     dsr = W_c'*dtanhInput;
150     ds_cur = dsr.*r(:, wordJ);
151     dsigInput_r = dsr.*s(:, wordJ-1).*r(:, wordJ).*(1-r(:, wordJ));
152     db_r = db_r + dsigInput_r;
153     dU_r = dU_r + dsigInput_r*x(:, wordJ)'; %could be accelerated by
avoiding add 0
154     dW_r = dW_r + dsigInput_r*s(:, wordJ-1)';
155     ds_cur = ds_cur + W_r'*dsigInput_r;
156
157     ds_cur = ds_cur + ds_cur_bk.*z(:, wordJ);
158     dz = ds_cur_bk.*(s(:, wordJ-1)-c(:, wordJ));
159     dsigInput_z = dz.*z(:, wordJ).*(1-z(:, wordJ));
160     db_z = db_z + dsigInput_z;

```



```

161     dU_z = dU_z + dsigInput_z*x(:,wordJ)'; %could be accelerated by
avoiding add 0
163     dW_z = dW_z + dsigInput_z*s(:,wordJ-1)';
ds_cur = ds_cur + W_z'*dsigInput_z;
165     end

166     % s_1
167     ds_cur = ds_cur + ds_single(:,1);

169     dtanhInput = (ds_cur.*(1-z(:,1)).*(1-c(:,1).*c(:,1)));
db_c = db_c + dtanhInput;
171     dU_c = dU_c + dtanhInput*x(:,1)'; %could be accelerated by avoiding
add 0
dW_c = dW_c + dtanhInput*(s_0.*r(:,1))';
173     dsr = W_c'*dtanhInput;
ds_0 = ds_0 + dsr.*r(:,1);
175     dsigInput_r = dsr.*s_0.*r(:,1).*(1-r(:,1));
db_r = db_r + dsigInput_r;
177     dU_r = dU_r + dsigInput_r*x(:,1)'; %could be accelerated by avoiding
add 0
dW_r = dW_r + dsigInput_r*s_0';
179     ds_0 = ds_0 + W_r'*dsigInput_r;

181     ds_0 = ds_0 + ds_cur.*z(:,1);
dz = ds_cur.*(s_0-c(:,1));
183     dsigInput_z = dz.*z(:,1).*(1-z(:,1));
db_z = db_z + dsigInput_z;
185     dU_z = dU_z + dsigInput_z*x(:,1)'; %could be accelerated by avoiding
add 0
dW_z = dW_z + dsigInput_z*s_0';
187     ds_0 = ds_0 + W_z'*dsigInput_z;
end

189 % A more direct view of backward propagate to calculate gradient using
% chain rule. (O(sentence_size^2) time)
% Instead of calculating how much contribution of derivative each step
has,
193 % here we calculate the gradient within every step.
function [dV, db_V, dU_z, dU_r, dU_c, dW_z, dW_r, dW_c, db_z, db_r, db_c,
ds_0] = ...
195 backward_direct(x, y, U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V,
b_V, s_0)
% forward propagate to get the intermediate and output results
197 [s, y_hat, L, z, r, c] = forward(x, y, U_z, U_r, U_c, W_z, W_r, W_c,
...
b_z, b_r, b_c, V, b_V, s_0);
199 % count sentence size
[~, sentence_size] = size(x);
201
202 % calculate gradient using chain rule
203 delta_y = y_hat - y;
db_V = sum(delta_y, 2);
205
206 dV = zeros(size(V));
207 for wordI = 1:sentence_size
dV = dV + delta_y(:,wordI)*s(:,wordI)';
209 end

211 ds_0 = zeros(size(s_0));
dU_c = zeros(size(U_c));
213 dU_r = zeros(size(U_r));
dU_z = zeros(size(U_z));
215 dW_c = zeros(size(W_c));
dW_r = zeros(size(W_r));
217 dW_z = zeros(size(W_z));

```

```

219 db_z = zeros(size(b_z));
220 db_r = zeros(size(b_r));
221 db_c = zeros(size(b_c));
222 ds_single = V'*delta_y;
223 % calculate the derivatives in each step and add them up
224 for wordI = 1:sentence_size
225     ds_cur = ds_single(:,wordI);
226     % since in each step t, the derivatives depends on s_0 - s_t,
227     % we need to trace back from t ot 0 each time
228     for wordJ = wordI:-1:2
229         ds_cur_bk = ds_cur;
230
231         dtanhInput = (ds_cur.*(1-z(:,wordJ)).*(1-c(:,wordJ)).*c(:,
wordJ)));
232         db_c = db_c + dtanhInput;
233         dU_c = dU_c + dtanhInput*x(:,wordJ)'; %could be accelerated
234 by avoiding add 0
235         dW_c = dW_c + dtanhInput*(s(:,wordJ-1).*r(:,wordJ))';
236         dsr = W_c'*dtanhInput;
237         ds_cur = dsr.*r(:,wordJ);
238         dsigInput_r = dsr.*s(:,wordJ-1).*r(:,wordJ).*(1-r(:,wordJ));
239         db_r = db_r + dsigInput_r;
240         dU_r = dU_r + dsigInput_r*x(:,wordJ)'; %could be accelerated
241 by avoiding add 0
242         dW_r = dW_r + dsigInput_r*s(:,wordJ-1)';
243         ds_cur = ds_cur + W_r'*dsigInput_r;
244
245         ds_cur = ds_cur + ds_cur_bk.*z(:,wordJ);
246         dz = ds_cur_bk.*(s(:,wordJ-1)-c(:,wordJ));
247         dsigInput_z = dz.*z(:,wordJ).*(1-z(:,wordJ));
248         db_z = db_z + dsigInput_z;
249         dU_z = dU_z + dsigInput_z*x(:,wordJ)'; %could be accelerated
250 by avoiding add 0
251         dW_z = dW_z + dsigInput_z*s(:,wordJ-1)';
252         ds_cur = ds_cur + W_z'*dsigInput_z;
253
254     end
255
256     % s_1
257     dtanhInput = (ds_cur.*(1-z(:,1)).*(1-c(:,1)).*c(:,1)));
258     db_c = db_c + dtanhInput;
259     dU_c = dU_c + dtanhInput*x(:,1)'; %could be accelerated by
260 avoiding add 0
261     dW_c = dW_c + dtanhInput*(s_0.*r(:,1))';
262     dsr = W_c'*dtanhInput;
263     ds_0 = ds_0 + dsr.*r(:,1);
264     dsigInput_r = dsr.*s_0.*r(:,1).*(1-r(:,1));
265     db_r = db_r + dsigInput_r;
266     dU_r = dU_r + dsigInput_r*x(:,1)'; %could be accelerated by
267 avoiding add 0
268     dW_r = dW_r + dsigInput_r*s_0';
269     ds_0 = ds_0 + W_r'*dsigInput_r;
270
271     ds_0 = ds_0 + ds_cur.*z(:,1);
272     dz = ds_cur.*(s_0-c(:,1));
273     dsigInput_z = dz.*z(:,1).*(1-z(:,1));
274     db_z = db_z + dsigInput_z;
275     dU_z = dU_z + dsigInput_z*x(:,1)'; %could be accelerated by
276 avoiding add 0
277     dW_z = dW_z + dsigInput_z*s_0';
278     ds_0 = ds_0 + W_z'*dsigInput_z;
279
280 end
281 end
282
283 % Sigmoid function for neural network
284 function val = sigmoid(x)

```

```

277     val = sigmf(x,[1 0]);
end

```

testBPTT_GRU.m

Less Important Functions

```

1 % Fake a training data set: generate only one sentence for training.
  %!!! Only for testing. Needs to be changed to read in training data from
  files.
3 function [x_t, y_t] = getTrainingData(vocabulary_size, sentence_size)
  assert(vocabulary_size > 2); % for start and end of sentence symbol
5  assert(sentence_size > 0);

7  % define start and end of sentence in the vocabulary
  SENTENCE_START = zeros(vocabulary_size, 1);
9  SENTENCE_START(1) = 1;
  SENTENCE_END = zeros(vocabulary_size, 1);
11 SENTENCE_END(2) = 1;

13 % generate sentence:
  x_t = zeros(vocabulary_size, sentence_size-1); % leave one slot for
  SENTENCE_START
15 for wordI = 1:sentence_size-1
  % generate a random word excludes start and end symbol
17   x_t(randi(vocabulary_size-2,1,1)+2, wordI) = 1;
  end
19 y_t = [x_t, SENTENCE_END]; % training output
  x_t = [SENTENCE_START, x_t]; % training input
21 end

23 % Use numerical differentiation to approximate the gradient of each
  % parameter and calculate the difference between these numerical results
25 % and our results calculated by applying chain rule.
function checkGradient_GRU(x, y, U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r,
  b_c, V, b_V, s_0, ...
27 dV, db_V, dU_z, dU_r, dU_c, dW_z, dW_r, dW_c, db_z, db_r, db_c, ds_0)
  % Here we use the centre difference formula:
29 % df(x)/dx = (f(x+h)-f(x-h)) / (2h)
  % It is a second order accurate method with error bounded by O(h^2)
31
  h = 1e-5;
33 % NOTE: h couldn't be too large or too small since large h will
  % introduce bigger truncation error and small h will introduce bigger
35 % roundoff error.

37 dV_numerical = zeros(size(dV));
  % Calculate partial derivative element by element
39 for rowI = 1:size(dV_numerical,1)
  for colI = 1:size(dV_numerical,2)
41     V_plus = V;
      V_plus(rowI, colI) = V_plus(rowI, colI) + h;
43     V_minus = V;
      V_minus(rowI, colI) = V_minus(rowI, colI) - h;
45     [~, ~, L_plus] = forward(x, y, ...
      U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V_plus, b_V,
47     s_0);
      [~, ~, L_minus] = forward(x, y, ...
      U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V_minus, b_V
49     , s_0);
      dV_numerical(rowI, colI) = (sum(L_plus) - sum(L_minus)) / 2 /
  h;
  end
51 end

```

```

53 display (sum(sum(abs(dV_numerical-dV)./(abs(dV_numerical)+h))), ...
    'dV relative error'); % prevent dividing by 0 by adding h
55 dU_c_numerical = zeros(size(dU_c));
56 for rowI = 1:size(dU_c_numerical,1)
57     for colI = 1:size(dU_c_numerical,2)
58         U_c_plus = U_c;
59         U_c_plus(rowI,colI) = U_c_plus(rowI,colI) + h;
60         U_c_minus = U_c;
61         U_c_minus(rowI,colI) = U_c_minus(rowI,colI) - h;
62         [~,~,L_plus] = forward(x,y,...
63             U_z, U_r, U_c_plus, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V,
        s_0);
64         [~,~,L_minus] = forward(x,y,...
65             U_z, U_r, U_c_minus, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V
        , s_0);
66         dU_c_numerical(rowI,colI) = (sum(L_plus) - sum(L_minus)) / 2
        / h;
67     end
68 end
69 display (sum(sum(abs(dU_c_numerical-dU_c)./(abs(dU_c_numerical)+h))),
    ...
    'dU_c relative error');
71 dW_c_numerical = zeros(size(dW_c));
72 for rowI = 1:size(dW_c_numerical,1)
73     for colI = 1:size(dW_c_numerical,2)
74         W_c_plus = W_c;
75         W_c_plus(rowI,colI) = W_c_plus(rowI,colI) + h;
76         W_c_minus = W_c;
77         W_c_minus(rowI,colI) = W_c_minus(rowI,colI) - h;
78         [~,~,L_plus] = forward(x,y,...
79             U_z, U_r, U_c, W_z, W_r, W_c_plus, b_z, b_r, b_c, V, b_V,
        s_0);
80         [~,~,L_minus] = forward(x,y,...
81             U_z, U_r, U_c, W_z, W_r, W_c_minus, b_z, b_r, b_c, V, b_V
        , s_0);
82         dW_c_numerical(rowI,colI) = (sum(L_plus) - sum(L_minus)) / 2
        / h;
83     end
84 end
85 display (sum(sum(abs(dW_c_numerical-dW_c)./(abs(dW_c_numerical)+h))),
    ...
    'dW_c relative error');
87 dU_r_numerical = zeros(size(dU_r));
88 for rowI = 1:size(dU_r_numerical,1)
89     for colI = 1:size(dU_r_numerical,2)
90         U_r_plus = U_r;
91         U_r_plus(rowI,colI) = U_r_plus(rowI,colI) + h;
92         U_r_minus = U_r;
93         U_r_minus(rowI,colI) = U_r_minus(rowI,colI) - h;
94         [~,~,L_plus] = forward(x,y,...
95             U_z, U_r_plus, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V,
        s_0);
96         [~,~,L_minus] = forward(x,y,...
97             U_z, U_r_minus, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V
        , s_0);
98         dU_r_numerical(rowI,colI) = (sum(L_plus) - sum(L_minus)) / 2
        / h;
99     end
100 end
101 display (sum(sum(abs(dU_r_numerical-dU_r)./(abs(dU_r_numerical)+h))),
    ...
    'dU_r relative error');

```

```

105 dW_r_numerical = zeros(size(dW_r));
107 for rowI = 1:size(dW_r_numerical,1)
109     for colI = 1:size(dW_r_numerical,2)
111         W_r_plus = W_r;
112         W_r_plus(rowI,colI) = W_r_plus(rowI,colI) + h;
113         W_r_minus = W_r;
114         W_r_minus(rowI,colI) = W_r_minus(rowI,colI) - h;
115         [~,~,L_plus] = forward(x,y,...
116             U_z, U_r, U_c, W_z, W_r_plus, W_c, b_z, b_r, b_c, V, b_V,
117             s_0);
118         [~,~,L_minus] = forward(x,y,...
119             U_z, U_r, U_c, W_z, W_r_minus, W_c, b_z, b_r, b_c, V, b_V,
120             s_0);
121         dW_r_numerical(rowI,colI) = (sum(L_plus) - sum(L_minus)) / 2
122         / h;
123     end
124 end
125 display(sum(sum(abs(dW_r_numerical-dW_r)./(abs(dW_r_numerical)+h)),
126     ...
127     'dW_r relative error'));
128
129 dU_z_numerical = zeros(size(dU_z));
130 for rowI = 1:size(dU_z_numerical,1)
131     for colI = 1:size(dU_z_numerical,2)
132         U_z_plus = U_z;
133         U_z_plus(rowI,colI) = U_z_plus(rowI,colI) + h;
134         U_z_minus = U_z;
135         U_z_minus(rowI,colI) = U_z_minus(rowI,colI) - h;
136         [~,~,L_plus] = forward(x,y,...
137             U_z_plus, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V,
138             s_0);
139         [~,~,L_minus] = forward(x,y,...
140             U_z_minus, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V,
141             s_0);
142         dU_z_numerical(rowI,colI) = (sum(L_plus) - sum(L_minus)) / 2
143         / h;
144     end
145 end
146 display(sum(sum(abs(dU_z_numerical-dU_z)./(abs(dU_z_numerical)+h)),
147     ...
148     'dU_z relative error'));
149
150 dW_z_numerical = zeros(size(dW_z));
151 for rowI = 1:size(dW_z_numerical,1)
152     for colI = 1:size(dW_z_numerical,2)
153         W_z_plus = W_z;
154         W_z_plus(rowI,colI) = W_z_plus(rowI,colI) + h;
155         W_z_minus = W_z;
156         W_z_minus(rowI,colI) = W_z_minus(rowI,colI) - h;
157         [~,~,L_plus] = forward(x,y,...
158             U_z, U_r, U_c, W_z_plus, W_r, W_c, b_z, b_r, b_c, V, b_V,
159             s_0);
160         [~,~,L_minus] = forward(x,y,...
161             U_z, U_r, U_c, W_z_minus, W_r, W_c, b_z, b_r, b_c, V, b_V,
162             s_0);
163         dW_z_numerical(rowI,colI) = (sum(L_plus) - sum(L_minus)) / 2
164         / h;
165     end
166 end
167 display(sum(sum(abs(dW_z_numerical-dW_z)./(abs(dW_z_numerical)+h)),
168     ...
169     'dW_z relative error'));
170
171 db_z_numerical = zeros(size(db_z));

```

```

159 for i = 1:length(db_z_numerical)
    b_z_plus = b_z;
    b_z_plus(i) = b_z_plus(i) + h;
161 b_z_minus = b_z;
    b_z_minus(i) = b_z_minus(i) - h;
163 [~, ~, L_plus] = forward(x, y, ...
        U_z, U_r, U_c, W_z, W_r, W_c, b_z_plus, b_r, b_c, V, b_V, s_0
    );
165 [~, ~, L_minus] = forward(x, y, ...
        U_z, U_r, U_c, W_z, W_r, W_c, b_z_minus, b_r, b_c, V, b_V,
s_0);
167 db_z_numerical(i) = (sum(L_plus) - sum(L_minus)) / 2 / h;
end
169 display(sum(abs(db_z_numerical-db_z)./(abs(db_z_numerical)+h)), ...
    'db_z relative error');
171
db_r_numerical = zeros(size(db_r));
173 for i = 1:length(db_r_numerical)
    b_r_plus = b_r;
175 b_r_plus(i) = b_r_plus(i) + h;
    b_r_minus = b_r;
177 b_r_minus(i) = b_r_minus(i) - h;
    [~, ~, L_plus] = forward(x, y, ...
179 U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r_plus, b_c, V, b_V, s_0
    );
    [~, ~, L_minus] = forward(x, y, ...
181 U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r_minus, b_c, V, b_V,
s_0);
    db_r_numerical(i) = (sum(L_plus) - sum(L_minus)) / 2 / h;
183 end
    display(sum(abs(db_r_numerical-db_r)./(abs(db_r_numerical)+h)), ...
185 'db_r relative error');
187
db_c_numerical = zeros(size(db_c));
189 for i = 1:length(db_c_numerical)
    b_c_plus = b_c;
    b_c_plus(i) = b_c_plus(i) + h;
191 b_c_minus = b_c;
    b_c_minus(i) = b_c_minus(i) - h;
193 [~, ~, L_plus] = forward(x, y, ...
        U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c_plus, V, b_V, s_0
    );
195 [~, ~, L_minus] = forward(x, y, ...
        U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c_minus, V, b_V,
s_0);
197 db_c_numerical(i) = (sum(L_plus) - sum(L_minus)) / 2 / h;
end
199 display(sum(abs(db_c_numerical-db_c)./(abs(db_c_numerical)+h)), ...
    'db_c relative error');
201
db_V_numerical = zeros(size(db_V));
203 for i = 1:length(db_V_numerical)
    b_V_plus = b_V;
205 b_V_plus(i) = b_V_plus(i) + h;
    b_V_minus = b_V;
207 b_V_minus(i) = b_V_minus(i) - h;
    [~, ~, L_plus] = forward(x, y, ...
209 U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V_plus, s_0
    );
    [~, ~, L_minus] = forward(x, y, ...
211 U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V_minus,
s_0);
    db_V_numerical(i) = (sum(L_plus) - sum(L_minus)) / 2 / h;
213 end
    display(sum(abs(db_V_numerical-db_V)./(abs(db_V_numerical)+h)), ...

```

```

215     'db_V relative error');
217     ds_0_numerical = zeros(size(ds_0));
219     for i = 1:length(ds_0_numerical)
221         s_0_plus = s_0;
221         s_0_plus(i) = s_0_plus(i) + h;
221         s_0_minus = s_0;
221         s_0_minus(i) = s_0_minus(i) - h;
223         [~, ~, L_plus] = forward(x, y, ...
223             U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V, s_0_plus
225     );
225         [~, ~, L_minus] = forward(x, y, ...
225             U_z, U_r, U_c, W_z, W_r, W_c, b_z, b_r, b_c, V, b_V,
225         s_0_minus);
227         ds_0_numerical(i) = (sum(L_plus) - sum(L_minus)) / 2 / h;
227     end
229     display(sum(abs(ds_0_numerical - ds_0) ./ (abs(ds_0_numerical) + h)), ...
229         'ds_0 relative error');
231 end

```

testBPTT_GRU.m