# Report on Project 3

Zhunxuan Wang, 13300180086
School of Mathematical Sciences

November 28, 2017

## 1   Recurrent Neural Network

### 1.1   Basic RNN Cell

RNN is an efficient model to deal with sequential information. For each node in the same layer, *recurrent* means they perform the same task for every element in the sequence, with the output being depended on the previous computations (memory) [3]
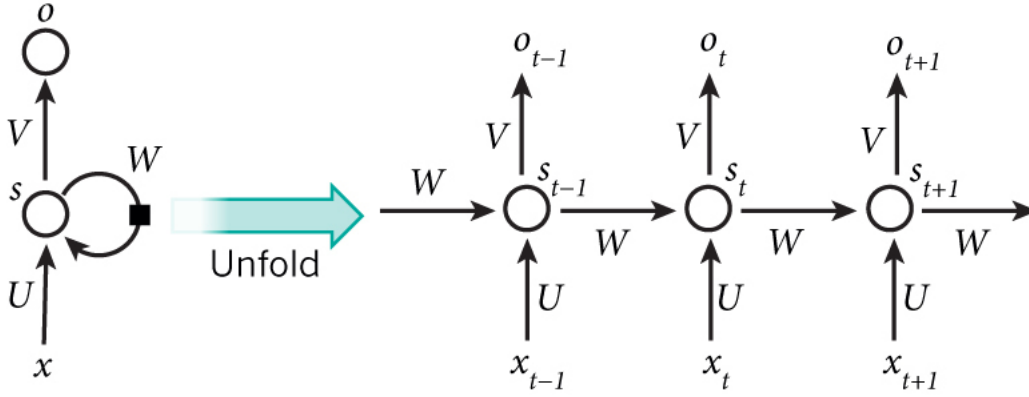


Figure 1: The RNN and its Unfolding

The above diagram shows the unfolding of basic RNN, of which each node shares the same parameters. The formulas dominating the calculation are expressed as follows

- $\mathbf{x}_t$ is the input vector (i.e. word embedding) at time step $t$.

- $\mathbf{s}_t$ is the hidden state (memory) at time step $t$, which is dominated by the previous hidden state and the input vector

$$\mathbf{s}_t = f\left(\mathbf{W}\mathbf{s}_{t-1} + \mathbf{U}\mathbf{x}_t\right)$$

  where $f$ is generally a non-linear function (e.g. tanh, `ReLu`)

- $\mathbf{o}_t$ is the output vector, which is dominated by the current hidden state

$$\mathbf{o}_t = g\left(\mathbf{V}\mathbf{s}_t\right)$$

  where $g$ depends on what we need (e.g. `softmax` for obtaining the probability of each words).

- $\mathbf{U}$, $\mathbf{V}$, $\mathbf{W}$ are the shared parameters, which are the parameters of a basic RNN cell.

## 1.2 Basic LSTM Cell

Long short-term memory (LSTM) [2] is a special kind of RNN, which is capable of remember long term dependencies. A basic LSTM cell is shown as follows
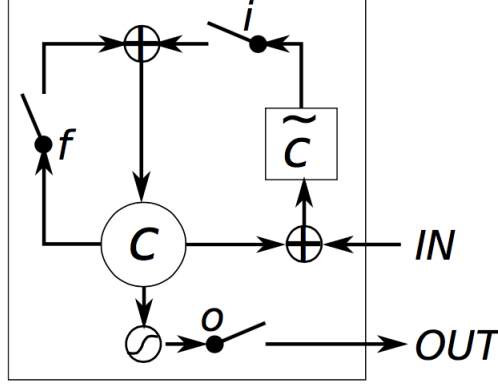


Figure 2: Basic Structure of an LSTM Cell

As the diagram shows, we introduce $\mathbf{i}_t$, $\mathbf{f}_t$, $\mathbf{o}_t$, and $\mathbf{c}_t$, $\tilde{\mathbf{c}}_t$ of which the definition and the formula are expressed as follows

- $\mathbf{i}_t$, $\mathbf{f}_t$, $\mathbf{o}_t$ are the input, forget and output gates respectively. They have the same form but different parameters

$$\mathbf{i}_t = \sigma\left(\mathbf{W}_i \mathbf{s}_{t-1} + \mathbf{U}_i \mathbf{x}_t\right)$$
$$\mathbf{f}_t = \sigma\left(\mathbf{W}_f \mathbf{s}_{t-1} + \mathbf{U}_f \mathbf{x}_t\right)$$
$$\mathbf{o}_t = \sigma\left(\mathbf{W}_o \mathbf{s}_{t-1} + \mathbf{U}_o \mathbf{x}_t\right)$$

  where $\mathbf{s}_t$, $\mathbf{x}_t$ are the hidden state and the input vector at time step $t$, $\mathbf{W}_{i,f,o}$ and $\mathbf{U}_{i,f,o}$ are the parameters of each gate.

- $\tilde{\mathbf{c}}_t$ is the new candidate value, which could be added to the state

$$\tilde{\mathbf{c}}_t = \tanh\left(\mathbf{W}_c \mathbf{s}_{t-1} + \mathbf{U}_c \mathbf{x}_t\right)$$

  of which the form is similar to the gates, but with a different activation function tanh.

- $\mathbf{c}_t$ is the internal memory, which is the sum of previous memory $\mathbf{c}_{t-1}$ multiplied by the forget gate $\mathbf{f}_t$, and the candidate value $\tilde{\mathbf{c}}_t$ multiplied by the input gate $\mathbf{i}_t$

$$\mathbf{c}_t = \mathbf{c}_{t-1} \circ \mathbf{f}_t + \tilde{\mathbf{c}}_t \circ \mathbf{i}_t$$

  where $\circ$ is the operator of elementwise muliplication.

- $\mathbf{s}_t$ is obtained by
$$\mathbf{s}_t = \tanh\left(\mathbf{c}_t\right) \circ \mathbf{o}_t.$$

Thereby the structure of a basic LSTM cell is established.

## 1.3  Basic GRU Cell

The structure behind a basic GRU (Gated Recurrent Unit) cell is quite similar to that of a basic LSTM cell
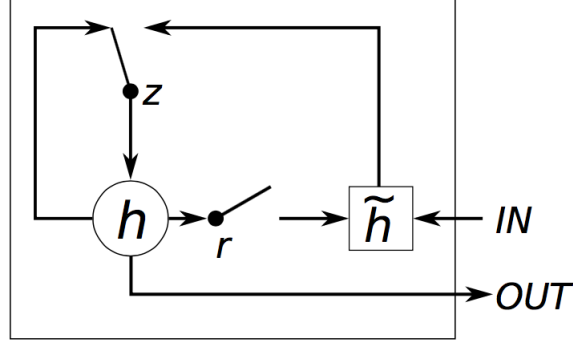


Figure 3: Basic Structure of a GRU Cell

As the diagram shows, we introduce $\mathbf{z}_t$, $\mathbf{r}_t$, and $\mathbf{h}_t$, $\tilde{\mathbf{h}}_t$ of which the definition and the formula are expressed as follows [1]

- $\mathbf{z}_t$, $\mathbf{r}_t$ are the update and reset gates respectively. They have the same form but different parameters

$$\mathbf{z}_t = \sigma\left(\mathbf{W}_{\mathrm{z}}\mathbf{s}_{t-1} + \mathbf{U}_{\mathrm{z}}\mathbf{x}_t\right)$$
$$\mathbf{r}_t = \sigma\left(\mathbf{W}_{\mathrm{r}}\mathbf{s}_{t-1} + \mathbf{U}_{\mathrm{r}}\mathbf{x}_t\right)$$

  where $\mathbf{s}_t$, $\mathbf{x}_t$ are the hidden state and the input vector at time step $t$, $\mathbf{W}_{\mathrm{z,r}}$ and $\mathbf{U}_{\mathrm{z,r}}$ are the parameters of each gate.

- $\tilde{\mathbf{h}}_t$ is the new candidate activation

$$\tilde{\mathbf{h}}_t = \tanh\left(\mathbf{W}_{\mathrm{h}}\left(\mathbf{s}_{t-1} \circ \mathbf{r}_t\right) + \mathbf{U}_{\mathrm{h}}\mathbf{x}_t\right).$$

  where $\circ$ is the operator of elementwise muliplication.

- $\mathbf{s}_t$ is obtained by

$$\mathbf{s}_t = \mathbf{s}_{t-1} \circ \mathbf{z}_t + \tilde{\mathbf{h}}_t \circ \left(\mathbf{1} - \mathbf{z}_t\right).$$

Thereby the structure of a basic GRU cell is established.

## 1.4  `Seq2seq` Generator

To generate poems using RNN model, we chose to adopt `seq2seq` model [4], of which the main idea of this model is to feed the current input by the previous output. We adopted 2-layer RNN model to generate poems
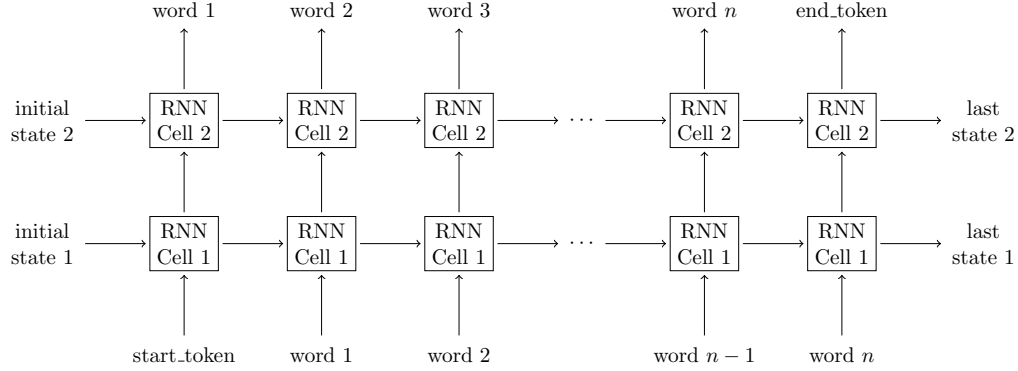


Figure 4: An Example of Poems Generating

As we observe from the graph above, the generating began with feeding input by a start token. After that, the input of each iteration would be fed by the previous output, and the incoming state would be the previous state. When it came to the last word, the generating would end with an end token.

It is worth mentioning that, if we specify a word to begin, the input after the start token will be the specific word.

# 2  Experimental Results

Setting the loss function the cross entropy of prediction and one-hot target, choosing LSTM as the RNN cell, the 2-layer RNN model was trained after 50 epochs, based on the dataset `poems.txt`.

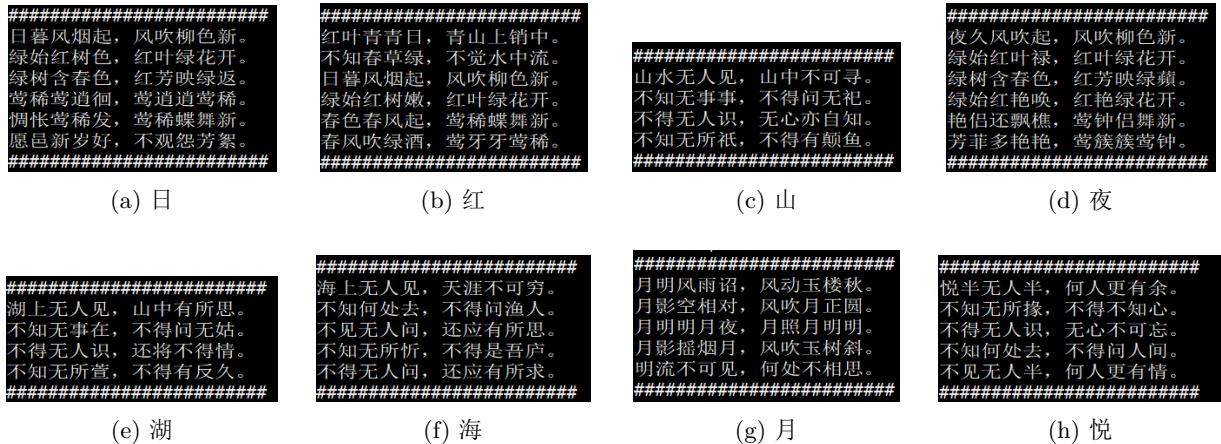we specified 8 words to start with, and the poems generated are shown as follows



Figure 5: LSTM-Generated Poems Based on `poems.txt`

As we observe from the poems above, they are all five-character poems, and the statement of the generated poems is relatively fluent. However the logic in the context did not make much sense, and too many words repeated.

4

Retraining the model based on the previous model and the new `tangshi.txt` dataset, a new set of poems were generated, as the figures above show


(a) 日


(b) 红


(c) 山


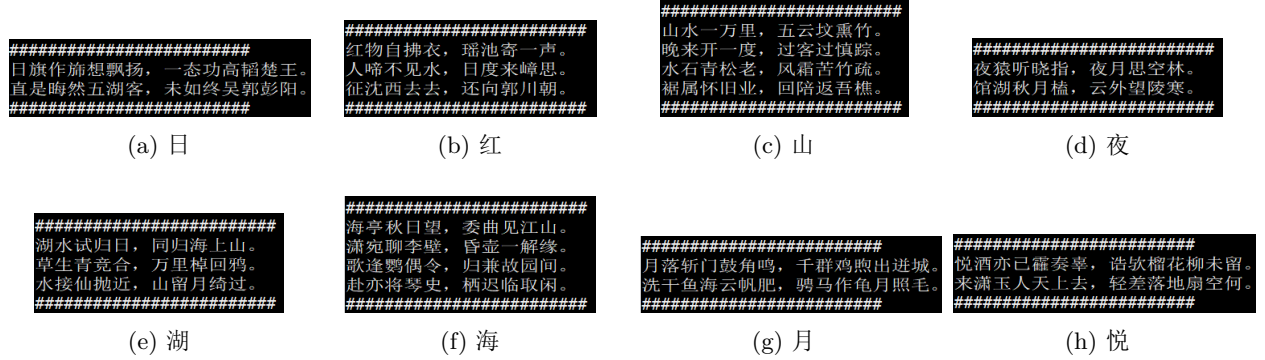(d) 夜


(e) 湖


(f) 海


(g) 月


(h) 悦

Figure 6: LSTM-Generated Poems Based on `poems.txt` and `tangshi.txt`

As we observe from the poems above, seven-character poems appeared, and the statement of the generated poems and the logic in the context is more fluent than before. Furthermore, the number of repeating words reduced a lot, and the efficiency on antithesis improved. In general, the model based on `poems.txt` and `tangshi.txt` is more efficient than that based on only `poems.txt`.

# 3 Comparison between basicRNN, LSTM and GRU

LSTM and GRU are gated mutants of RNN structure. They all perform well on poem generator model training, and better than basic RNN structure, with reference to the total train loss mean in the last epoch (when reaching the convergence state) $L$ (as the following table shows, only on `poems.txt`)

|   | basicRNN | LSTM | GRU |
|---|---|---|---|
| $L$ | 5.819501 | 4.082350 | 4.222888 |

From the table we can infer that the performance by LSTM and GRU are close (LSTM are a little better than GRU), each of which performs better than basic RNN (on `poems.txt`).

The poems generated by basic RNN and GRU are in the appendix A. As we can obtain from the poems, the poems generated by basic RNN are very short and monotonous, and the poems generated by GRU make more sense that by basic RNN, but they lack more logic in context and have more repeating words than LSTM-generated poems.

# References

[1] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.

[2] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[3] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.

[4] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.

# A   Poems by basicRNN and GRU



(a) 日     (b) 红     (c) 山     (d) 夜

(e) 湖     (f) 海     (g) 月     (h) 悦

Figure 7: basicRNN-Generated Poems Based on `poems.txt`



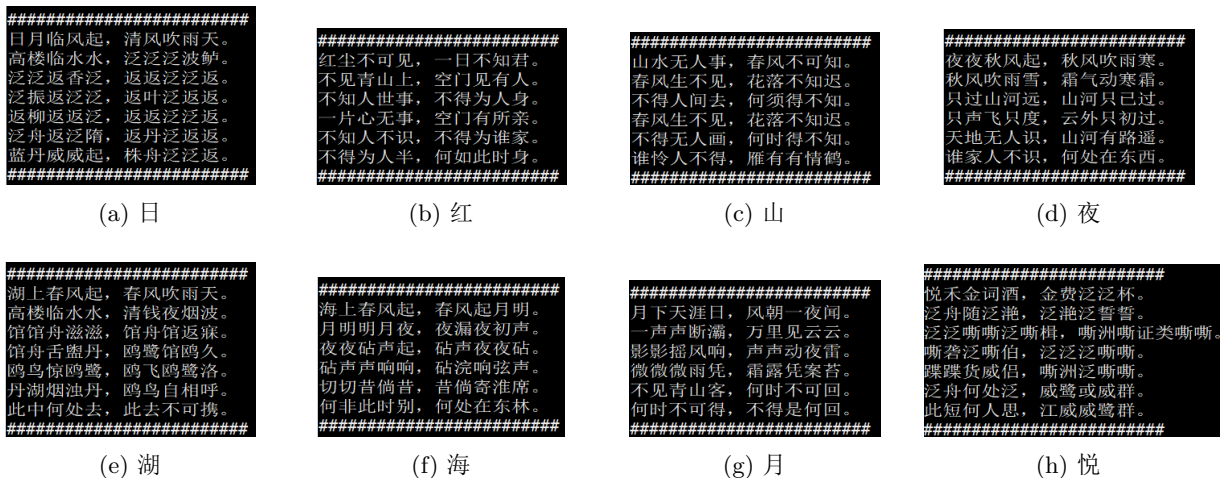(a) 日     (b) 红     (c) 山     (d) 夜

(e) 湖     (f) 海     (g) 月     (h) 悦

Figure 8: GRU-Generated Poems Based on `poems.txt`