

Layerwise Interweaving Convolutional LSTM

Tiehang Duan and Sargur N. Srihari

Department of Computer Science and Engineering
The State University of New York at Buffalo
Buffalo, NY 14260, United States
tiehangd@buffalo.edu, srihari@cedar.buffalo.edu

Abstract. A deep network structure is formed with LSTM layer and convolutional layer interweaves with each other. The Layerwise Interweaving Convolutional LSTM(LIC-LSTM) enhanced the feature extraction ability of LSTM stack and is capable for versatile sequential data modeling. Its unique network structure allows it to extract higher level features with sequential information involved. Experiment results show the model achieves higher accuracy and shoulders lower perplexity on sequential data modeling tasks compared with state of art LSTM models.

1 Introduction

LSTM is a specific Recurrent Neural Network(RNN) which uses explicit gates to control the magnitude of input, cell state and hidden state at each time step and it has the unique property of maintaining long term dependencies in the memory cell [6]. It is playing an important role in sequential learning and has been applied to multiple fields of natural language processing including sentiment analysis [8], semantic analysis [7] and neural machine translation [1] [5] etc. The feature extraction of LSTM can be enhanced with windowed convolution operation, which is used widely in computer vision and image processing [4]. Researchers have contributed innovative and pioneering work for possible combinations of LSTM and CNN. Y. Kim et al. built the CNN-LSTM character level model for NLP problems [3]; S. Vosoughi et al. combined CNN and LSTM model to encode tweets into vectors for semantic and sentiment analysis [7]; and J. Donahue et al. used similar structure for visual recognition and description [2].

Previous works adopt the network structure with convolutional layer(or several convolutional layers) at the bottom and LSTM layers on the top, and feature extraction process is completed before it is fed into the LSTM stack. In this paper, we introduced dense convolutional layer between all adjacent LSTM layers in the network and forms the Layerwise Interweaving Convolutional LSTM(LIC-LSTM).

2 Network Structure and Mechanism

The network structure is shown in Fig. 1. In the net, Convolution layer and LSTM layer interweaves with each other. The features extracted in the filters

are combined and fed into the next LSTM layer, and the encoded output from the LSTM layer is convoluted to extract higher level features containing sequential information. The evolving of feature from low level to higher level is more homogeneous. While max pooling or mean pooling is often applied in CNN to reduce feature dimensionality, no pooling is introduced in our network based on two considerations: (1) pooling disrupts the sequential order [10], (2) LSTM is capable to deal with long term sequences and shrinking of time steps is not necessary.

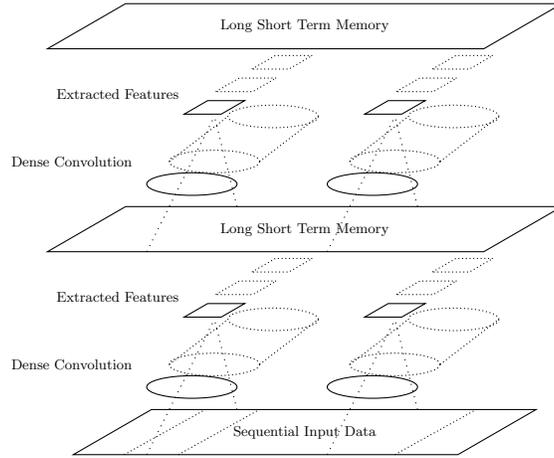


Fig. 1: Structure of Layerwise Interweaving Convolutional LSTM(LIC-LSTM)

Valid convolution is used for each convolutional layer in the forward pass. Meaning of labels is summarized in Table 1.

The extracted feature after the convolutional layer is

$$f_i^n = \sum_{j=1}^C \left(\sum_{k=1}^D W[f]_{jk}^n \times I_{(i+C-j)(D+1-k)} \right) \quad (1)$$

where n is the n th filter and i denotes the i th time step.

The extracted features from N different filters are combined together to form the new input to LSTM layer, in which the following transformations are performed

$$i_{gate} = \text{sigmoid}(W_{ih}H_{s-1} + W_{ix}I_s + b_i) \quad (2)$$

$$o_{gate} = \text{sigmoid}(W_{oh}H_{s-1} + W_{ox}I_s + b_o) \quad (3)$$

$$f_{gate} = \text{sigmoid}(W_{fh}H_{s-1} + W_{fx}I_s + b_f) \quad (4)$$

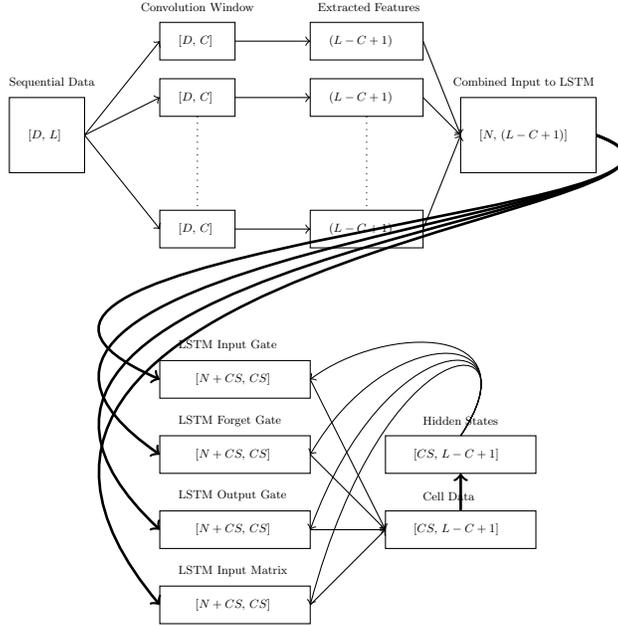


Fig. 2: Workflow in one layer of LIC-LSTM, dimension of feature and data is illustrated in the brackets.

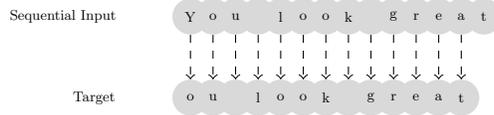


Fig. 3: Illustration of Character Prediction

$$C_s = f_{gate} \times C_{s-1} + i_{gate} \times \tanh(W_{in}[I_s; H_{s-1}]) \quad (5)$$

$$H_s = o_{gate} \times \tanh(C_s) \quad (6)$$

where i_{gate} , f_{gate} , o_{gate} are referred to as input, forget and output gates. Weights of the LSTM are W_{jk} with $j \in \{i, f, o, in\}$ and $k \in \{h, x\}$. C_s denotes the memory cell and H_s denotes the hidden state. The evolving of feature for one layer of LIC-LSTM is shown in Fig. 2.

3 Network Training and Performance

We reveal the detailed training process of the network with character prediction task. For the majority of related works, the training of CNN stack and LSTM stack are processed separately. This approach, while easy to implement and fine tune, can not guarantee overall optimization for the whole combined network.

Table 1: Summary of Label Meaning

Label	Meaning
D	Character Representation Dimension
L	Number of Time Steps
I	Input Data
W[f]	Filter Weight
C	Convolution Layer Window Size
N	Number of Filters
CS	Dimension of Memory Cell and Hidden State in LSTM

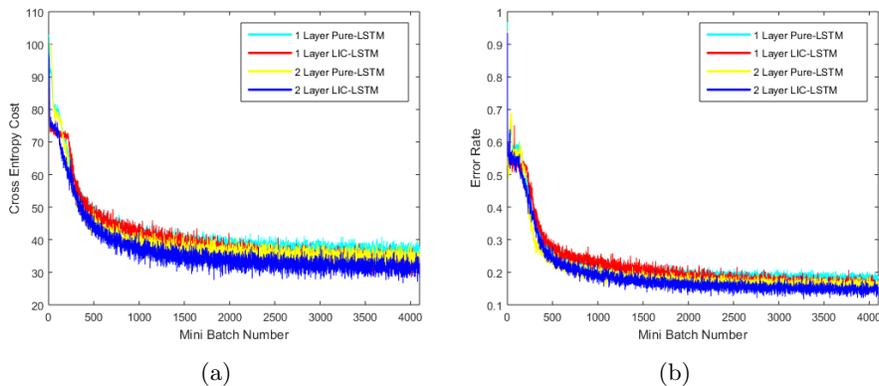


Fig. 4: Performance Comparison of Different Models. (a) Cross Entropy Cost, (b) Error Rate.

In LIC-LSTM, the convolutional layers and LSTM layers interweave with each other, and in the forward backward training phase, the gradients are propagated through the whole network simultaneously and concurrently. Such approach is designed to achieve overall optimization in the whole network and better performance. The drawback is requisite for more elegant implementation and careful fine tuning.

We used 24,000 English phrases for training and 8,000 English phrases for testing. The dataset is from LightNet toolbox introduced in [9]. The labels are created by moving the training phrase one character left. An illustration of the task is shown in Fig. 3. A total of 67 different characters are included and each character is encoded into one hot vector. The number of time steps is set to 34 and the original phrases are either truncated or concatenated with zeros. SGD is used for the convolutional layer optimization and Adam is used for the LSTM layer. As some of the characters like spaces are much more frequent than other characters, we adopt soft error rate in our experiment. The prediction is classified as correct if the label appears in the top K predictions ($K = 5$ in our implementation). The momentum is set to be 0.5 initially and then increased

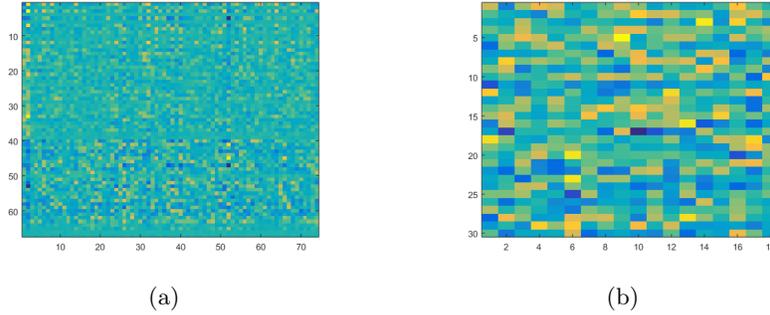


Fig. 5: Visualization of Filter Weight. (a) First Convolution Layer, (b) Second Convolution Layer.

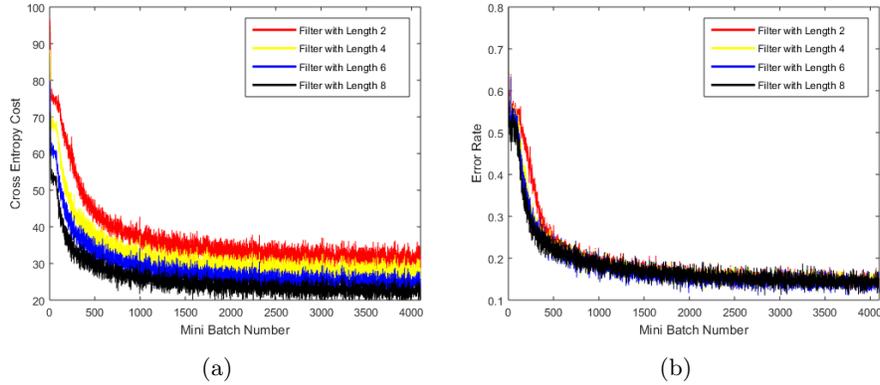


Fig. 6: Performance comparison with different filter length. (a) 2 layer LIC-LSTM cross entropy cost, (b) 2 layer LIC-LSTM error rate.

to 0.9 after 20 mini batches. The initial low momentum allows the model to rapidly reach functioning region and the following high momentum avoids the model being trapped in local optima. We found high momentum(0.9) is necessary for LIC-LSTM with 2 layers to achieve satisfying result. The model training performance is shown in Fig. 4. For the same number of layers, the LIC-LSTM model is achieving higher accuracy than pure LSTM model. And the performance of deep LSTM model lies between shallow LIC-LSTM model and deep LIC-LSTM model. The deep LIC-LSTM model achieves 16.2% error rate on testing dataset compared to 17.8% achieved by deep LSTM model. The LIC-LSTM model is enduring less cross entropy cost with larger filter size as shown in Fig. 6, while error rate remains stable, which shows the model performance is robust on varying filter size.

4 Conclusion

In this paper, we formed LIC-LSTM with convolutional layer and LSTM layer interweaving with each other. LIC-LSTM combines the feature extracting ability of CNN and sequential data modeling of LSTM homogeneously. Experiment results show the model achieves higher accuracy and shoulders lower cross entropy cost on sequential data modeling tasks compared with state of art LSTM models. The LIC-LSTM can serve as encoders in encoder-decoder model with natural language processing and neural machine translation tasks.

References

1. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. CoRR abs/1409.0473 (2014), <http://arxiv.org/abs/1409.0473>
2. Donahue, J., Hendricks, L.A., Rohrbach, M., Venugopalan, S., Guadarrama, S., Saenko, K., Darrell, T.: Long-term recurrent convolutional networks for visual recognition and description. IEEE Transactions on Pattern Analysis and Machine Intelligence PP(99), 1–1 (2016)
3. Kim, Y., Jernite, Y., Sontag, D., Rush, A.M.: Character-aware neural language models. CoRR abs/1508.06615 (2015), <http://arxiv.org/abs/1508.06615>
4. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Pereira, F., Burges, C.J.C., Bottou, L., Weinberger, K.Q. (eds.) Advances in Neural Information Processing Systems 25, pp. 1097–1105. Curran Associates, Inc. (2012), <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
5. Luong, M., Pham, H., Manning, C.D.: Effective approaches to attention-based neural machine translation. CoRR abs/1508.04025 (2015), <http://arxiv.org/abs/1508.04025>
6. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N.D., Weinberger, K.Q. (eds.) Advances in Neural Information Processing Systems 27, pp. 3104–3112. Curran Associates, Inc. (2014), <http://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf>
7. Vosoughi, S., Vijayaraghavan, P., Roy, D.: Tweet2vec: Learning tweet embeddings using character-level CNN-LSTM encoder-decoder. CoRR abs/1607.07514 (2016), <http://arxiv.org/abs/1607.07514>
8. Wang, J., Yu, L.C., Lai, K.R., Zhang, X.: Dimensional sentiment analysis using a regional cnn-lstm model. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, pp. 225–230. Berlin, Germany (2016)
9. Ye, C., Zhao, C., Yang, Y., Fermüller, C., Aloimonos, Y.: Lightnet: A versatile, standalone matlab-based environment for deep learning. CoRR abs/1605.02766 (2016), <http://arxiv.org/abs/1605.02766>
10. Zhou, C., Sun, C., Liu, Z., Lau, F.C.M.: A C-LSTM neural network for text classification. CoRR abs/1511.08630 (2015), <http://arxiv.org/abs/1511.08630>