

ENGINEERING SCIENCE

RNN WITH ADDITIONAL CONSTANT MEMORY FOR IMAGE CAPTION GENERATION TASK

Poghosyan Aghasi,
Sarukhanyan Hakob

Armenia, Yerevan, Institute for Informatics and Automation Problems of NAS RA

ARTICLE INFO

Received 24 May 2017

Accepted 05 June 2017

Published 05 July 2017

KEYWORDS

Supervised learning,
deep learning,
image caption generation,
RNN,
LSTM.

ABSTRACT

Analyze and generation of automated captions for images is one of the most common problems in artificial intelligence. Existing works use LSTM (Long Short-Term Memory) as recurrent neural network cell to solve this task. After training their deep neural models can generate image caption. But there is an issue, the next predicted word of the caption depends mainly on the last predicted word, rather than image content. In this paper, we present model that can automatically generate image description and is based on a recurrent neural network with modified LSTM cell that has an additional gate responsible for image features. This modification results in generation of more accurate captions. We have trained and tested our model on MSCOCO image dataset by using only images and their captions.

© 2017 The Authors.

1. INTRODUCTION

Image caption generation via automated system is a complex task aimed to support visually impaired people in better understanding of the content of images on the web. This approach can also have big impact on search engines and in robotics. This task is significantly harder than well-studied image classification [1] or object recognition.

Image caption must contain not only image object names, but their properties, relations, and actions. Moreover, the generated caption must be expressed through natural language like English. Which means, that already pre-trained neural language model needs an additional visual information to generate image caption.

Number of works have been published approaching this problem. Some of them [2] [3] [4] offer to combine existing image object detection and sentence generation systems. But there is the more efficient solution [5] that offers a joint model that takes an image and generates the caption, which describes image adequately. Last achievements in statistical machine translation were actively used in image caption generation tasks. The reason for this mainly is the proven achievement of greater results when using a powerful sequential model trained by maximizing the probability of the correct translation for the input sentence. These models [6] [7] [8] are based on Recurrent Neural Networks (RNNs). The model encodes variable

length input into fixed length vector representation. This representation enables conversion of the input sentence into the target sentence or the input image into the target image caption. The last model was being trained to maximize $p(S|I)$ likelihood to generate the target sequence of words $S = \{S_1, S_2, \dots\}$ for an input image I , where each word S_t comes from a given dictionary, that describes the image adequately.

Public image datasets lack in detailed descriptions and variety of image scenes. These limitations cause RNN to predict subsequent word of the sentence by ignoring the image scene or objects. Thus, next word prediction mainly depends on the previous word.

MODEL

Models based on machine translation that can generate image descriptions actively use a recurrent neural network. We will maximize the probability of the correct caption for the given image.

$$\theta^* = \arg \max_{\theta} \sum_{(I,S)} \log p(S|I; \theta) \quad (1)$$

In formulation (1) θ represents the parameters of our model and S is its correct caption for the given image I . If we have a sentence $S = \{S_0, S_1, \dots, S_N\}$ with the length of N , then we can apply the chain rule to calculate joint probability (2) over S_0, S_1, \dots, S_N ,

$$\log p(S|I; \theta) = \sum_{t=0}^N \log p(S_t|I, S_0, \dots, S_{t-1}; \theta) \quad (2)$$

where (I, S) is a training example pair. While training, we optimize the sum of the log probabilities for the whole training set using stochastic gradient descent [9]. $p(S_t|I, S_0, \dots, S_{t-1}; \theta)$ probability will correspond to the t step (iteration) of Recurrent Neural Network (RNN) based model. The variable number of words that are conditioned upon, up to $t - 1$ is expressed by a fixed length hidden state or memory h_t . After every iteration for the new input, x_t memory will be updated (3) by using a non-linear function f .

$$h_{t+1} = f(h_t, x_t) \quad (3)$$

For f we use a Long-Short Term Memory (LSTM), which has shown state-of-the-art performance on sequence generation tasks, such as translation or image caption generation. Model consists of the feed forward deep convolutional neural network (CNN) that feeds RNN.

One of the best Convolutional Neural Networks (CNN) is *Google Inception* [10], that has been widely used in object classification and object detection tasks. Furthermore, there are works [11] that have done CNN transfer learning for object classification for such tasks as scene classification.

There are high-level features that describe image semantic content like objects, their properties and relations in CNN [12]. In this work, we will select *Mixed_7c* layer from *Google Inception* and append *average pooling* layer which will have 2048-dimensional output for image description. Also, we will append *fully connected* neural layer with N_e neurons, which will convert 2048-dimensional vector into N_e dimensional vector. N_e is a words embedding vector's dimensionality [13]. The output vector x_{-1} of fully connected layer will be the first feed vector for RNN.

$$i_t = \sigma(W_{ix}x_t + W_{im}m_{t-1})$$

$$c_t = f_t \odot c_{t-1} + i_t \odot h(W_{cx}x_t + W_{cm}m_{t-1}) \quad (1)$$

$$f_t = \sigma(W_{fx}x_t + W_{fm}m_{t-1})$$

$$m_t = o_t \odot c_t \quad (2)$$

$$o_t = \sigma(W_{ox}x_t + W_{om}m_{t-1})$$

$$p_{t+1} = \text{softmax}(W_{pm} * m_t) \quad (3)$$

In (4-9) equations i_t, o_t, f_t are input, output and forget gates correspondingly, c_t is a cell memory in step t and m_t is an output of the LSTM for step (iteration) t . $W_{ix}, W_{im}, W_{fx}, W_{fm}, W_{ox}, W_{om}, W_{cx},$ are trainable parameters (variables) of the LSTM. \odot represents the product with a gate

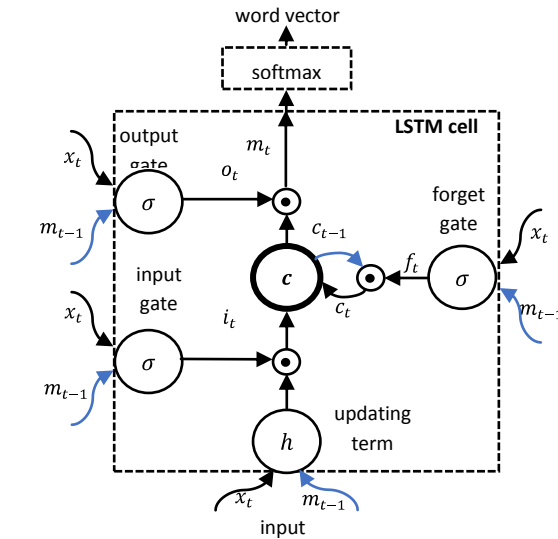


Fig. 1. LSTM: the memory block contains a cell c which is controlled by three gates

Long Short-Term Memory (LSTM) is an RNN cell. It helps in solving RNN training time problems like vanishing and exploding gradients [14], which is a significant problem for RNNs. LSTM is commonly used in machine translation, sequence generation and image description generation tasks. Work [5] uses recurrent neural network with an LSTM cell to generate image caption.

Constructionally LSTM is a memory cell c encoding knowledge at every iteration of what inputs have been seen up to this iteration. Later this knowledge is used for subsequent word generation (8, 9). Behavior of the cell is controlled by three gates – *input gate*, *output gate* and *forget gate*. Each gate is a vector of real number elements ranging from 0 to 1. In particular (see Fig. 1), forget gate is responsible for controlling whether to forget the cell's old value, input gate controls the permission for reading a new input value and finally output gate controls the permission to output the new value from the cell. This is done by multiplying the given gate with corresponding value (4, 5, 6). The definition of the LSTM are as follows:

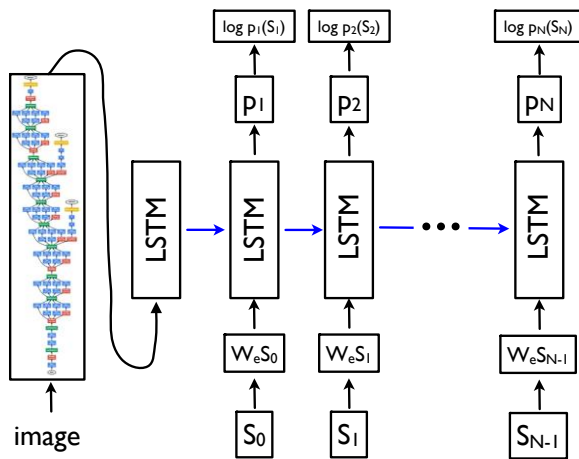


Fig. 2. LSTM model combined with a CNN image embedder (as defined in) and words embedding

The LSTM model is trained to predict the probability for the next word of an image caption after it has observed all previous words in the captions and image features. For easier training LSTM is represented in unrolled form (see **Ошибка! Источник ссылки не найден.**), which is a copy of the LSTM memory for the image and each word of the sentence. Also all LSTMs share the same parameters. Thus, x_{-1} is the first input for the first LSTM. Initial state of the LSTM is c_{-1} zero-filled memory. For the next LSTMs, inputs correspond to the word embedded vectors. Also, all recurrent connections are converted into feed-forward connections.

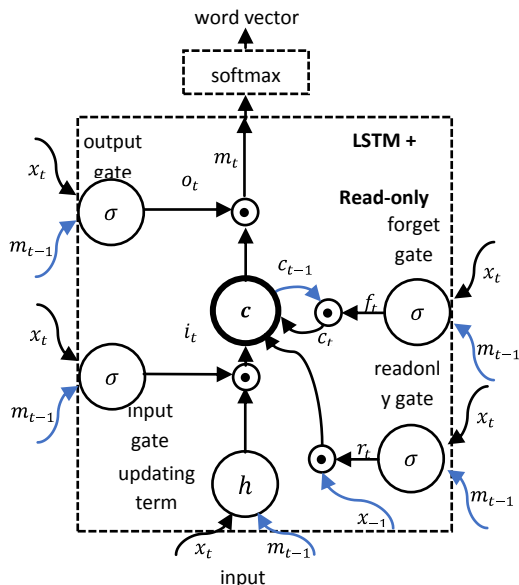


Fig. 3. LSTM: the memory block contains a cell c which is controlled by four gates (additional read-only memory)

For the input image I and the image's true caption $S = \{S_0, S_1, \dots, S_N\}$, the unrolling procedure is:

$$x_{-1} = CNN(I), x_t = W_e S_t, p_{t+1} = LSTM(x_t)$$

where each word S_t is the row of square identity $N_d \times N_d$ matrix at corresponding index, where N_d is the dictionary size. Also, S_0 is a special start word and S_N is a special stop word which indicates the start and the end of the sentence.

Note that both the image and the words are mapped to the same space. Vision CNN's last fully connected layer maps the image content to the embedding space. Also words are embedded by words embedding W_e , where W_e is a trainable parameter with $N_e \times N_d$ dimensionality.

Loss [5] is the sum of the negative log likelihood of the correct word at each step:

$$L(I, S) = - \sum_{t=1}^N \log p(S_t)$$

After training the model by minimizing loss with gradient descent, we will have all the parameters of the LSTM, the top layer of the image embedder CNN and word embedding W_e .

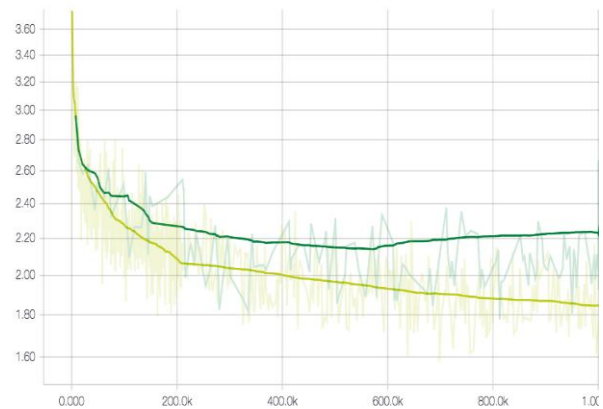


Fig. 3. Train loss function, LSTM – green, LSTM with read-only memory – yellow

After one million training iterations on MSCOCO [15] image dataset, we got loss function graphic (see green graphic on **Ошибка! Источник ссылки не найден.**). Experiments have shown that the LSTM's next prediction mainly depends on the previous word, that is why the LSTM generates words that are not associated with an input image. We will add an additional unit to the LSTM, that will help to create the new LSTM state that depends on the image content. This additional unit will be a new gate whose

value will be appended in the state calculation (see **Ошибка! Источник ссылки не найден.**).

$$r_t = \sigma(W_{rx}x_t + W_{rm}m_{t-1}) \quad (7)$$

$$c_t = f_t \odot c_{t-1} + r_t \odot x_{-1} + i_t \odot h(W_{cx}x_t + W_{cm}m_{t-1}) \quad (8)$$

In (14) r_t is the read-only gate with W_{rx} and W_{rm} additional trainable parameters. New

state c_t is calculated as shown in (15). After retraining with the new gate, we can see the new training loss (see yellow graphic on **Ошибка! Источник ссылки не найден.**).

We inference by using BeemSearch and already trained parameters to generate image caption as presented in work [5]. Some examples of inference are presented bellow (see **Ошибка! Источник ссылки не найден.**).





			
<u>LSTM</u>	<u>LSTM</u>	<u>LSTM</u>	<u>LSTM</u>
a) a group of people sitting around a table eating food. b) a group of people sitting around a wooden table. c) a group of people sitting around a table with food.	a) a woman sitting on a beach with an umbrella. b) a woman sitting on a beach chair holding an umbrella. c) a woman sitting on a beach with an umbrella	a) a man and a woman sitting on a bench . b) a couple of people that are sitting down d) a man and a woman sitting on a bench with their cell phones.	a) a woman sitting at a table with a glass of wine. b) a woman holding a glass of wine in her hand. c) a woman sitting at a table with a glass of wine
<u>LSTM + Read-only cell</u>	<u>LSTM + Read-only cell</u>	<u>LSTM + Read-only cell</u>	<u>LSTM + Read-only cell</u>
a) a group of people sitting around a picnic table. b) a group of people sitting around a table . c) a group of people sitting around a table eating food.	a) a woman sitting on the beach under an umbrella. b) a woman sitting on the beach with an umbrella. c) a woman sitting on the beach under an umbrella	a) a man and a woman sitting on a bus. b) a man and a woman are sitting on a bus. d) a man and a woman are sitting on a train.	a) a man is holding a glass of wine. b) a man is holding a glass of wine c) a man holding a glass of wine in front of a glass.

Fig. 4. Image caption generated by LSTM and LSTM with Read-only Unit models

CONCLUSION

In this work various existing automated caption generation systems have been analyzed in order to create new Recurrent Neural Network (RNN) with Long Short-Term Memory (LSTM) cell and a Read-only Unit has been developed. An additional unit has been added to work [5] that increases the model accuracy. Two models, one with the LSTM and other with the LSTM and Read-

only Unit have been trained on the same MSCOCO image train dataset. The best (the loss is minimum) middle loss values are 2.15 for LSTM and 1.85 for LSTM with Read-only Unit. MSCOCO image test dataset has been used for testing. Loss values for LSTM and LSTM with Read-only Unit model test are 2.05 and 1.90 accordingly. These metrics have shown that the new RNN model can generate image caption more accurately.

REFERENCES

1. O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," arXiv:1409.0575, 2014.

2. Farhadi, M. Hejrati, M. A. Sadeghi, P. Young, C. Rashtchian, J. Hockenmaier and D. Forsyth, "Every picture tells a story: Generating sentences from images," in *ECCV*, 2010.
3. G. Kulkarni, V. Premraj, S. Dhar, S. Li, Y. Choi, A. C. Berg and T. L. Berg, "Baby talk: Understanding and generating simple image descriptions," *CVPR*, 2011.
4. Karpathy and L. Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions," *IEEE*, vol. 39, no. 4, pp. 664 - 676, 2015.
5. O. Vinyals, A. Toshev, S. Bengio and D. Erhan, "Show and Tell: Lessons learned from the 2015 MSCOCO Image Captioning Challenge," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, no. 99, July 2016.
6. D. Bahdanau, K. Cho and Y. Bengio, "Neural machine translation by jointly learning to align and translate," arXiv:1409.0473, 2014.
7. K. Cho, B. van Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation.," *EMNLP*, 2014.
8. Sutskever, O. Vinyals and Q. Le, "Sequencetosequence learning with neural networks.," *NIPS*, 2014.
9. L. Bottou, "Large-scale machine learning with stochastic gradient descent.," in *Physica-Verlag HD*, 2010.
10. C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens and Z. Wojna, "Rethinking the inception architecture for computer vision," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2818--2826, 2016.
11. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng and T. D. DeCaf, "A deep convolutional activation feature for generic visual recognition," in *ICML*, 2014.
12. Poghosyan and H. Sarukhanyan, "Image Visual Similarity Based on High Level Features of Convolutional Neural Networks," *Mathematical Problems of Computer Science*, vol. 45, pp. 138--142, 2016.
13. T. Mikolov, K. Chen, G. Corrado and J. Dean, "Efficient Estimation of Word Representations in Vector Space," in *Proceedings of Workshop at ICLR*, 2013.
14. S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, 1997.
15. T. Lin, M. Maire, S. Belongie, J. Hays and P. Perona, "Microsoft coco: Common objects in context.," in *European Conference on Computer Vision*, 2014.