

RNAprotein binding sites and motifs prediction

LUO LU

HAN CHENYU

11155107885

11551xxxxx

1155107885@link.cuhk.edu.hk

jane@smith.com

December 2, 2020

Abstract

RNA Binding proteins(RBPs) are a class of proteins in eukaryotes and play important roles in many biological process, including gene splicing. Unfortunately, reseaches about RBPs, especially the binding preferences, are far from sufficient, so the prediction performance of traditional bioinformatics algorithms is not satisfactory due to the deficiency of prior information. Thanks to the development of high-throughput technology, vast experiment results are generated in recent few years, providing the precondition of using data-driven approaching like machine learning and CNNs. In this report, we design a model which consists of two CNNs, train it for each RBPs in the dataset, test it performance and compare it with other similar solutions finally. We also discuss difference between our model and other solutions in both the structure and performance, and try to give some potential improvements.

I. INTRODUCTION

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

II. BACKGROUND

i. RBPs

ii. CNN

CNN is a kind of supervised learning that use multiple convolutional layer to extract hidden feature from input data. CNN is widely used in image recognition and classification algorithm. A typical CNN consists of convolutional, pooling and full connected layers. Convolutional layers are used to extract features, pooling layers are used to decrease the dimensions of data and full connected layers can choose which features are more relative to which class.

III. METHODS

The key of our method is that a specific RBP can only recognize one class of RNA sequences which have similar subsequence in somewhere of the rna, called motif. Therefore, our goal is to recognize the similar subsequence, motif, from the given RNA sequence.

One method is based on the experiment results. As long as we get some motifs from experiments, we can use traditional bioinformatics tools like k-mers to calculate the similarity of input sequence and known motifs. However, the high-throughput technologies are time-intensive and expensive [Hafner, Markus, et al, 2009], and results are not satisfactory, so the number of accurately known motif is small compared to the real situation. And other drawback is that this method may not be able to make a accurate prediction if the corresponding motif hasn't been found.

The second method is to use data-driven approaching like machine learning. Our project is based on this method. In our method, we don't need the knowledge of motifs, and the network will try to find the features we need. CNN, as a powerful feature extracting tool,

has potential to find similar structure from giving data, and make a more accurate prediction compared to the first method.

IV. IMPLEMENTATION

There are lots of reseaches showing that deep learning can perform much better than the previous method [Pan, Xiaoyong, et al, 2018]. There are several network has promising performance, including CNNs and LSTMs. In our project, we decide to use two kinds of CNNs, global CNN and local CNN, to make the prediction, and here are the implementation details:

- **Global CNNs:** The purpose of this CNN is to extract useful feature from RNA sequence. In this CNN, we will feed the whole RNA sequence to the network, so that the network will find features no matter where the subsequence locate. Since the CNN can only receive fixed length input, we need to read and preprocess the data before training and test models(Details about the data preprocess can be found in the third points). For the network, we set 2 convolutional layers, 2 pooling layers and 2 dense layers. The kernel size of the first layer is [10, 4], because as Deepbind suggested, the best kernel length is about 1.5 times of the average motifs' length which 7 [Alipanahi, Babak, et al, 2015]. The kernel size of the second layer is [10, 1] due to the same reason.
- **Local CNNs:** The purpose of the CNN is to find the relationship of consecutive features. In this network, we will feed the multi-channel tensor to the network, so that the relationship between different channel can make difference. Since the raw data contain only one channel, we also need to do the preprocessing. For this network, we set 3 convolutional layers, 3 pooling layer and 2 dense lay-

ers. This networks is more complicated because the dimensions of input data are greater, and the complexity of data is higher.

- **Data preprocessing:** We downloaded the RBP-24 dataset from the website of GraphProt <http://www.bioinf.uni-freiburg.de/Software/GraphProt>. The origin data type in this dataset is nucleotide sequence. In order to feed the data to CNNs, we need to transform the string into one-hot tensor. After transforming, every sequence can be seen as a graph with width of 4. For global CNNs, we will find the max sequence length among data of each RBPs, and then padding the rest to that length by adding $[[0.25],[0.25],[0.25],[0.25]]$ to the tail of each one-hot matrix. For local CNNs, we will divide the sequence into fixed length subsequence with fixed shifting length, and then append the later subsequence to the channel of previous subsequence, until the channels are full or there is no subsequence left. In our project, the window size has been set to 101, shifting size is 20 and the channel size is 7.
- **Models ensembling:** Model ensembling is common because it's hard to get a high performance predictor by using only one model. In our situation, we also meet this problem that the prediction is not accurate if we only use one of the CNN models, because both the sequence character and the features independence are important for binding. Therefore, we train those 2 CNN separately, and use both in the prediction phase, by just average the output probabilities(probabilities of positive and negative). In our experiment, this method can increase the accuracy.
- **Other details:** Due to the limitation of experiments, the number of positive data and negative data is unbalanced. In order to ensure the bal-

ance between positive input and negative input, we use over-sampling strategy, which means the sample in the small size class can be used much more times in training. This method can significantly increase the training efficient [Ando, Shin and Huang, 2017].

V. RESULTS

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

$$e = mc^2 \quad (1)$$

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a,

egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

VI. DISCUSSION

i. Unsatisfactory prediction result for dataset *Baltz2012*

From the result part showed, the prediction for first 3 RBPs is not accurate enough. We discussed about why these performance are different, and give some potential reason:

- **Overfitting:** Actually, we have found this problem in the early phase of our project, and the reason we gave at that time is overfitting because the training accuracy was not bad but the testing accuracy is super low, even lower than 50 in that time. High training accuracy and low testing accuracy are the character of overfitting. However, we have cut down the number of layers and neurons per layer since finding that situation, and the final version of our networks is not complex. However, the problem have been improved but not really solved. Therefore, overfitting may be parts of the reason, but we don't think it's the main reason.
- **Small size of datasets:** We also found that the datasize for those 3 RBPs is significantly smaller than others, especially the first dataset. We natually guessed that the size is too small to represent all the features to recognize the special structures.
- **Quality of Experiment result:** As we mentioned before, the high-throughput technology is an expensive and time-intensive new technology, so the quality of a experiment may not as reliable as some mature technology. We found that not only our project, others also meet the similar problem. Maybe there are too many false result in those datasets, which cannot be distinguished by our project.

ii. Subsection Two

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

REFERENCES

- [Hafner, Markus, et al, 2009] Hafner, Markus, et al. "Transcriptome-wide identification of RNA-binding protein and microRNA target sites by PAR-CLIP." *Cell* 141.1 (2010): 129-141.
- [Pan, Xiaoyong, et al, 2018] Pan, Xiaoyong, et al. "Prediction of RNA-protein sequence and structure binding preferences using deep convolutional and recurrent neural networks." *BMC genomics* 19.1 (2018): 511.
- [Alipanahi, Babak, et al, 2015] Alipanahi, Babak, et al. "Predicting the sequence specificities of DNA-and RNA-binding proteins by deep learning." *Nature biotechnology* 33.8 (2015): 831-838.
- [Ando, Shin and Huang, 2017] Ando, Shin, and Chun Yuan Huang. "Deep over-sampling framework for classifying imbalanced data." *Joint European Conference on Machine Learning and Knowledge*

Discovery in Databases. Springer, Cham,
2017.