

ON SECONDARY STRUCTURE ANALYSIS BY USING FORMAL GRAMMARS AND ARTIFICIAL NEURAL NETWORKS

Polina Lunina^(1,2), Semyon Grigorev^(1,2)

(1) Saint Petersburg State University, 7/9 Universitetskaya nab., St. Petersburg, 199034, Russia
lunina_polina@mail.ru, s.v.grigoriev@spbu.ru

(2) JetBrains Research, Primorskiy prospekt 68-70, Building 1, St. Petersburg 197374, Russia
semyon.grigorev@jetbrains.com

Keywords: DNN, CNN, Machine Learning, Secondary Structure, Genomic Sequences, Formal Grammars, Parsing.

Abstract. Recently a way to combine formal grammars and artificial neural networks for biological sequences processing was proposed. This approach utilizes a grammar for encoding the primitive features of RNA secondary structure (in contrast to the classical way, when probabilistic grammars are used for modeling the secondary structure of the whole sequence). Subsequently, these features are extracted by parsing algorithm and processed by neural network to perform some sort of classification. In this work we provide further development of the proposed approach. Also we show that it is possible to create a model that handles original sequences and does not require parsing in practical usage. The idea here is to perform two-staged learning: first, training a neural network that classifies parsed data, and second, extending it with a number of input layers that transform the nucleotide sequence into parsing result. We evaluate the proposed improvements on some tRNA classification tasks and show that these improvements are applicable while using our approach and demonstrate high practical performance.

1 Introduction

Developing of effective computational methods for genomic sequences analysis is an open problem in bioinformatics. The existing algorithms for sequences classification and subsequences detection adopt different concepts and approaches but the fundamental idea here is that secondary structure of genomic sequences contains important information about the biological functions of organisms. There are different ways of secondary structure formal description such as probabilistic grammars, covariance models and Hidden Markovs Models [1, 2, 3].

The common problem while dealing with real-world data is a possible presence of different mutations, noises and random variations which requires some sort of probability estimation while modeling the secondary structure. Probabilistic grammars and covariance models provide such functionality along with good expressive possibilities and long-distance connections handling, and they are successfully used in some tools, such as [4], but building and training accurate grammar or model for predicting the whole secondary structure involves some theoretical and practical difficulties [?].

In [5] an approach for biological sequences processing using the combination of formal grammars and artificial neural networks is proposed. The key idea is to use an ordinary context-free grammar to describe only the key secondary structure features and leave the probabilistic analysis to neural network which takes parsing-provided data as an input and performs some sort of classification. Neural networks are a common way to

process noisy data and find complex structural patterns, moreover, the efficiency of neural networks for genetic data processing have already been shown in some works [6, 7]. The applicability of the proposed approach for some real-world tasks was demonstrated and in the present work we provide some improvements and modifications in the context of transport RNA (tRNA) classification tasks.

2 Proposed solution

In this work we describe some new ideas that may provide more computational possibilities and improve the trained model accuracy, compared to the previous work. The first idea is to explore different formats of the parsing results representation and develop corresponding neural network architectures. The second idea is to minimize the use of parsing in the context of our solution.

While using the proposed approach on some genetic sequences classification task, the first step is to describe the main structural elements of the RNA secondary structure (stems and loops) by context-free grammar and extract them by means of parsing. Our solution is not dependent on a specific parsing algorithm, but here a matrix-based version [8] is used due to its high practical performance and the possibilities in use of parallel computing. The result of a parsing algorithm for the input string w and the fixed grammar non-terminal N can be presented as an upper-triangular boolean matrix M_N , where $M_N[i, j] = 1$, iff the substring $w[i, j - 1]$ is derivable from N . We use such matrices as an input to neural network that is supposed to perform classification of some kind by detecting sufficient features and finding patterns in their appearance. Therefore, we need to transform these boolean matrices to some data structures accepted by neural network. Presently, we came up with two possible ways. The first one is to drop out the bottom left triangle, vectorize it row by row and transform it to the byte vector. This approach reduces the input size, but it requires the equal length of the input sequences, therefore we propose to either cut sequences or add some special symbol for each of them till the definite length. The second way is to represent the matrix as an image: the false bits of matrix as white pixels and the true bits as black ones. This approach makes it possible to process sequences with different length since the images could be easily transformed to a constant size.

The final step is to process the parsing-provided data by an artificial neural network constructed and trained for a specific task. The neural network architecture is unique for each problem, but during the experimental research we worked out some common concepts. For vectorized data we use dense layers because data locality is broken during vectorization and dropout layers with batch normalization to stabilize learning. For image data we use a small number of convolutional layers, then linearization and go to the same architecture as for vectorized data, because convolution layers are suited for features extraction, but in our case it is already done by parsing. In the previous work we performed experiments only on vectorized data and in this work we provide an evaluation on both data formats and compare the results.

The bottleneck of our solution is parsing and the main problem here is following: in practical usage of trained model we need to parse the input sequence which is quite time-consuming while working with huge biological databases. To solve this problem we propose to use two-staged learning. Firstly, we train a neural network on the parsed data that performs classification according to a given problem. After that, we extend this neural network by a number of input layers that take the initial nucleotide sequence as an inputs and convert it to the parsing result. So, we build a model that handles sequences and requires parsing only for training the model it is based on. This way we can remove the parsing step from the practical usage of trained model and also improve its accuracy without the additional data generation. In the next section we provide the results of this modification usage.

3 Experiments

We evaluate the proposed approach with described above modifications on two tRNA sequences analysis tasks. The first one is a classification of tRNA into two classes: eukaryotes and prokaryotes. And the second one is a classification into four classes: archaea, bacteria, plants and fungi. For these experiments we use sequences from tRNA databases [9, 10]. For both classification tasks we took the equal amount of samples for each class and generated vectors and images datasets by parsing tool on the same sequences dataset using grammar presented in our previous work. After that, we trained neural networks that take parsed data (vectors or images) and perform classification for our tasks. Then, for both vector- and image-based models we created the extended neural network which contains two blocks: the first one takes initial tRNA sequence as an input and transforms it to the parsing result by a number of dense layers with batch normalization and the second one copies the sequence of layers from the base model and uses its weights while training. The example of such neural network architecture for two classes classification based on vectorized data is presented in figure 1, where the right rectangle is the original model that classifies vectors and the left rectangle is the extension that transforms sequence to vector. For image-based approach we use the same architecture, except in that case we removed the convolutional layer from the extended model, thus, at the junction of the blocks the first layer corresponds to linearized image.

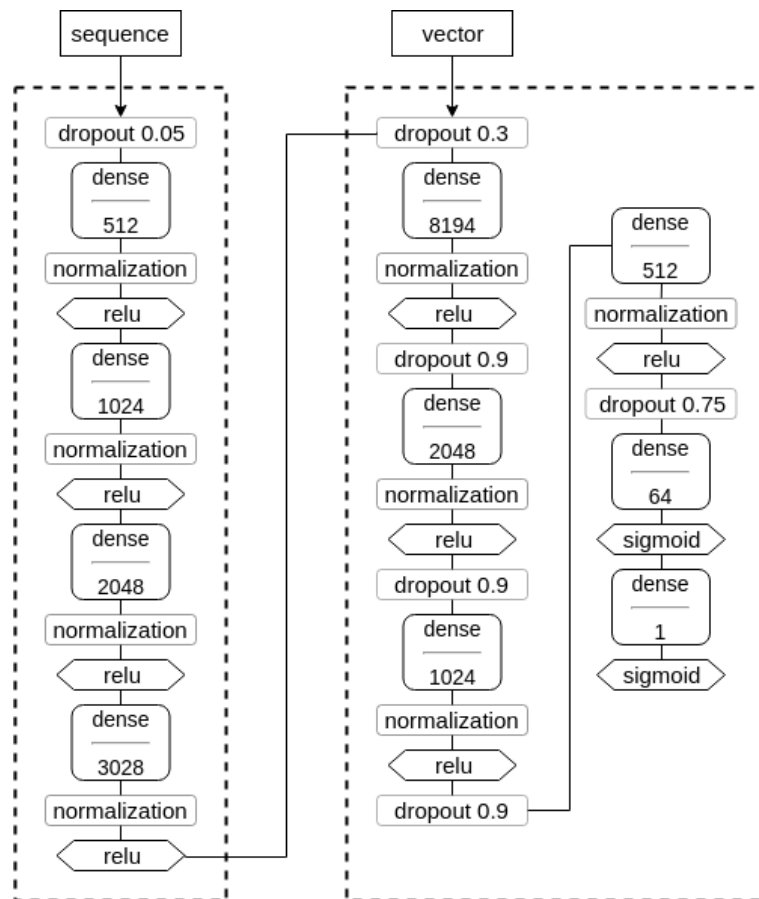


Figure 1: Neural network architecture

The trained models for two classes (EP) and for four classes (ABFP) classification tasks were evaluated by classical machine learning metrics: accuracy, precision and recall. Results on test dataset for each problem by accuracy metrics are presented in the table 1, where base model is a model which handles parsing result (image or vector re-

spectively) and extended model handles tRNA sequences and extends the corresponding base model.

Table 1: Models test results by accuracy metrics

Classifier	EP		ABFP	
Approach	Vector-based	Image-based	Vector-based	Image-based
Base model accuracy	94.1%	96.2%	86.7%	93.3%
Extended model accuracy	97.5%	97.8%	96.2%	95.7%
Total samples (train:valid:test)	20000:5000:10000		8000:1000:3000	

Test results for extended models for both classifiers on the same datasets as in table 1 by precision and recall metrics are presented in the table 2.

Table 2: Models test results by precision and recall metrics for each class

Classifier	Class	Vector-based approach		Image-based approach	
		precision	recall	precision	recall
EP	prokaryotic	95.8%	99.4%	96.2%	99.4%
	eukaryotic	99.4%	95.6%	99.4%	99.5%
ABFP	archaeal	91.1%	99.2%	91.6%	98.5%
	bacterial	96.6%	95.1%	95.2%	95.5%
	fungi	98.5%	94.9%	97.5%	94.3%
	plant	99.4%	95.7%	99.2%	94.7%

The results show that our approach is applicable to RNA classification tasks and both vector- and image-based models can be used along with dense and convolutional layers in neural networks architectures. The differences in results are insignificant, because this is basically the same data and neural network is able to process it in a similar way, so, the format of parsed data representation can be chosen based on a particular problem or some data processing specifics. Moreover, the idea of extended model that handles sequences is proved to be applicable in practice and it demonstrates even higher performance than the original parsing-based model, as it can be seen in the table 1.

4 Conclusion

We describe modifications of the proposed approach for biological sequences analysis using the combination of formal grammars and neural networks. We showed that it is possible to represent parsing result as an image and use convolutional layers while processing it with neural network. Also we developed a solution that removes the parsing step from the trained model use and allows to test models on the original RNA sequences. We demonstrated the applicability of the proposed modifications on real-world problems. Source code and documentation are published at GitHub: <https://github.com/LuninaPolina/SecondaryStructureAnalyzer>.

Acknowledgments

The research was supported by the Russian Science Foundation grant 18-11-00100 and a grant from JetBrains Research.

References

- [1] R. Durbin, S. R. Eddy, A. Krogh, and G. Mitchison, *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge university press, 1998.

- [2] R. D. Dowell and S. R. Eddy, "Evaluation of several lightweight stochastic context-free grammars for rna secondary structure prediction," *BMC bioinformatics*, vol. 5, no. 1, p. 71, 2004.
- [3] B. Knudsen and J. Hein, "Rna secondary structure prediction using stochastic context-free grammars and evolutionary history," *Bioinformatics (Oxford, England)*, vol. 15, no. 6, pp. 446–454, 1999.
- [4] E. P. Nawrocki and S. R. Eddy, "Infernal 1.1: 100-fold faster RNA homology searches," *Bioinformatics*, vol. 29, pp. 2933–2935, Nov 2013.
- [5] S. Grigorev and P. Lunina, "The composition of dense neural networks and formal grammars for secondary structure analysis,"
- [6] D. Sherman, "Humidor: Microbial community classification of the 16s gene by training cigar strings with convolutional neural networks," 2017.
- [7] S. Higashi, M. Hungria, and M. Brunetto, "Bacteria classification based on 16s ribosomal gene using artificial neural networks," in *Proceedings of the 8th WSEAS International Conference on Computational intelligence, man-machine systems and cybernetics*, pp. 86–91, 2009.
- [8] R. Azimov and S. Grigorev, "Context-free path querying by matrix multiplication," in *Proceedings of the 1st ACM SIGMOD Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA)*, GRADES-NDA '18, (New York, NY, USA), pp. 5:1–5:10, ACM, 2018.
- [9] "Genomic tRNA Database." <http://gtrnadb.ucsc.edu/>. Last accessed 05.06.2019.
- [10] "tRNADB-CE." <http://trna.ie.niigata-u.ac.jp/cgi-bin/trnadb/index.cgi>. Last accessed 05.06.2019.