



Broadband Long-Term Spectrum Prediction Based on Trend Based SAX

Han Zhang, Lu Sun, and Yun Lin^(✉)

College of Information and Communication, Harbin Engineering University, Harbin 150001,
People's Republic of China
liny@hrbeu.edu.cn

Abstract. With the development of communication technology and the growth of equipment, spectrum prediction technology has received more and more attention because of its wide application in spectrum resource management. However, due to the high burst of spectrum usage, there are still some difficulties in spectrum prediction. This paper proposes a new algorithmic framework (TrSAX-seq2seq) for the difficult problem of broadband and long-term prediction in the spectrum prediction problem. In this paper, a trend based symbolic aggregate approximation (TrSAX) method is used to reduce the dimension and represent the historical spectrum observation data, and then perform hierarchical clustering on the symbol sequence after dimension reduction to achieve the purpose of dividing the wide frequency band into multiple narrow frequency bands. Then, we use the LSTM network of seq2seq architecture to predict the spectrum occupation. We validate our method on a real spectrum monitoring dataset. The experimental results show that the method proposed in this paper can effectively improve the prediction accuracy compared with other methods.

Keywords: Broadband spectrum prediction · Long term prediction · Trend based symbolic aggregate approximation · Seq2seq

1 Introduction

The rapid development of communication technology and the sharp increase in the number of communication terminal equipment make the limited electromagnetic spectrum resources more scarce. Therefore, researchers pay more and more attention to the research of spectrum prediction technology, because the spectrum of a specific area can be dynamically allocated in hours or tens of hours based on the prediction of the future state of the spectrum [1]. However, the use of many spectrums is dynamic and highly bursty, and the spectrum usage rules of different frequency bands are very different because of their different services and users [2]. At the same time, we always want to predict as far into the future as possible so that dynamic spectrum allocation can be done earlier. These problems bring enormous challenges to spectrum prediction and lead us to consider new algorithmic frameworks to achieve long-term spectrum prediction in broadband scenarios.

At present, scholars have proposed many methods for spectrum prediction. In the traditional method, the ARMA algorithm is used more [3]. Among machine learning based methods, algorithms such as Support Vector Machine (SVM) [4], Long Short Term Memory (LSTM) [5], Convolutional Long Short Term Memory (ConvLSTM) [6], CNN + LSTM [14] are also applied to spectrum prediction. For the ISM frequency band with high burstiness, some people propose a clustering method based on statistical features to cut the spectrum and predict based on LSTM respectively [5]. Most of the above methods are used to predict the short-term busy and idle state of the spectrum [7]. The wide-band and long-term prediction problems mentioned above are still not well solved.

The problem of spectrum prediction is essentially a time series prediction problem. An effective idea is to cluster large sequences into multiple smaller sequences to improve prediction accuracy. But considering the long dimension of spectrum monitoring data, we must consider a dimensionality reduction method to ensure the effect of clustering. On the problem of long-term prediction of spectrum, we believe that the seq2seq architecture is more potential than the traditional LSTM network. Based on the above discussion, this paper uses a trend based symbolic aggregate approximation (TrSAX) method combined with a seq2seq network for spectrum prediction, and tests it on a real open source dataset. The main contributions of this paper are as follows:

1. This paper uses a trend based symbolic aggregate approximation method to represent and reduce the dimension of electromagnetic spectrum data. Then, multiple prediction models are constructed based on hierarchical clustering, and the wide-band spectrum prediction problem is transformed into multiple narrow-band prediction problems.
2. Apply the seq2seq network architecture to the spectrum prediction problem to improve the prediction performance in long-term prediction.
3. Validate the proposed method using real open-source spectrum monitoring data and demonstrate its superior performance.

The rest of this paper is outlined below. Section 2 introduces the trend based symbolic aggregate approximation method. Section 3 introduces the seq2seq architecture for spectrum prediction. Section 4 is the experimental part of this paper. Section 5 is the conclusion.

2 Trend Based Symbolic Aggregate Approximation Method

For spectrum monitoring data, due to its large data volume and high dimension, the effect of traditional data mining and clustering methods is affected. In recent years, academia has proposed many representation methods that have outstanding performance in time series, which are characterized by the ability to achieve dimensionality reduction without losing important features [8]. In this paper, a relatively new trend based symbolic aggregate approximation (TrSAX) method is used to reduce the dimension of spectrum monitoring data, and the numerical value and changing trend characteristics of the data are preserved [9]. The broadband electromagnetic spectrum data is then divided into multiple narrow frequency bands using hierarchical clustering.

2.1 Monitoring Data Autocorrelation Analysis

There have been many research results showing that there is widespread correlation between different channels [10]. The correlation of channels can be divided into positive correlation and negative correlation. This paper defines positive correlation as when one channel is busy (idle), the other channel is also in a busy (idle) state. A negative correlation is the opposite. In this paper, the channel correlation factor (CCF) is used to represent the degree of correlation between channels, and its calculation method is shown in formula (1). Where ρ is the correlation coefficient of the two channels, which is the value of CCF; M is the total number of channels, C_i^m is the state of channel i at time m , $I\{A\}$ is the index function, when A is true, $I\{A\} = 1$, otherwise $I\{A\} = 0$. The numerator of formula (1) is the same number of two channel states, and the denominator is the total number of channel states within the monitoring time.

$$\rho = \frac{\sum_{m=1}^M I\{C_i^m = C_j^m\}}{\sum_{m=1}^M I\{C_i^m = C_j^m\} + \sum_{m=1}^M I\{C_i^m \neq C_j^m\}} \quad (1)$$

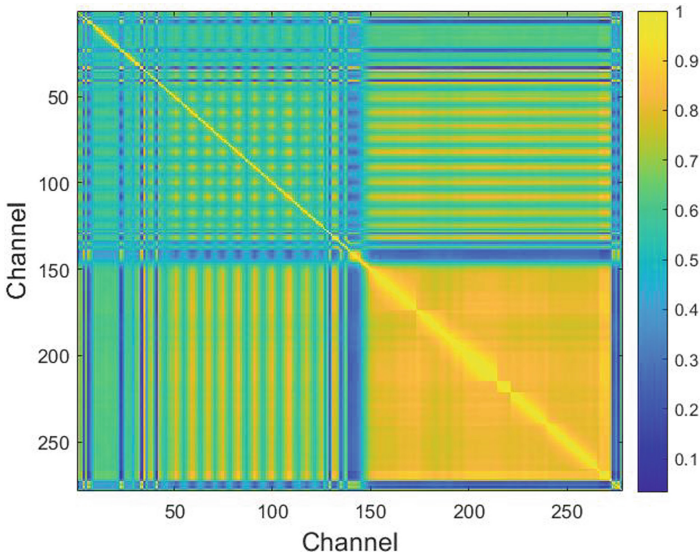


Fig. 1. Correlation between different channels.

Figure 1 is the correlation diagram of all channels in the experimental data in Chapter 4 of this paper. It can be seen that the correlation between channels is ubiquitous, and the CCF value between most channels is above 0.4. The CCF value between some channels is close to 1, which indicates that the evolution laws of these channels are highly correlated. Therefore, it is a feasible method to study a clustering method to find channels with similar laws and train prediction models separately.

2.2 Trend Based Symbolic Aggregate Approximation (TrSAX)

Based on the piecewise aggregation approximation (PAA) method, scholars further developed the symbolic aggregation approximation (SAX) and the trend based symbolic aggregation approximation (TrSAX) method used in this paper.

The TrSAX method assumes that the original time series data approximately obey the normal distribution, divides the original data into several time periods, and uses two symbols to represent the mean and slope of the time periods respectively, so as to achieve the purpose of dimension reduction. The specific process is as follows: first, the original time series data needs to be normalized to data with standard deviation of 1 and average value of 0. Then the original time series is divided into several time periods by using PAA method, and the average value of each time period is obtained. Then, the average value of each segment is compared with the breakpoint of the equal probability partition space under the Gaussian curve, and then mapped to the corresponding characters. Let the length of the time series $X = \{x_1, \dots, x_n\}$ be n . Represent it by a vector \bar{X} of length N , the i th element of \bar{X} is defined as:

$$\bar{X}_i = \frac{N}{n} \sum_{j=(n(i-1)+1)/N}^{ni/N} x_j \tag{2}$$

In this way, the symbolic representation of the original time series data is realized, and the fast and effective dimensionality reduction of the time series is realized. After dimensionality reduction, the overall change of time series can still be reflected. Figure 2 is a schematic diagram of a symbol mapping method.

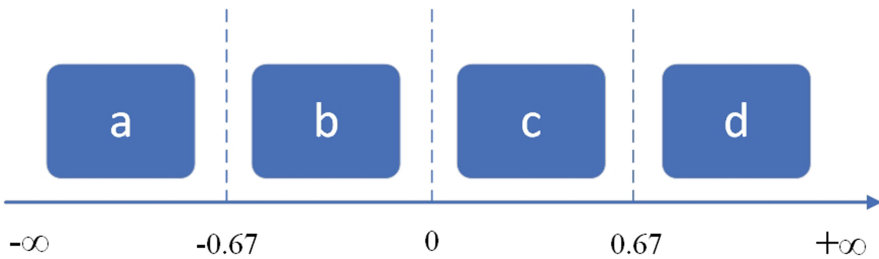


Fig. 2. Schematic diagram of the TrSAX symbol mapping method with word size 4.

The trsax method further adds the trend change information to the symbol coding. The method additionally calculates the slope between each time period. The main improvement of the method is to use two symbols in each time period. Where the first symbol is the segment average value in the original Sax method, and the second symbol is the segment slope value. In this article, lower case letters represent the segment average value, and upper case letters represent the segment slope value. The schematic diagram of segment slope mapping is shown in Fig. 3.

The segment slope value k is estimated by the least squares method:

$$k = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} \tag{3}$$

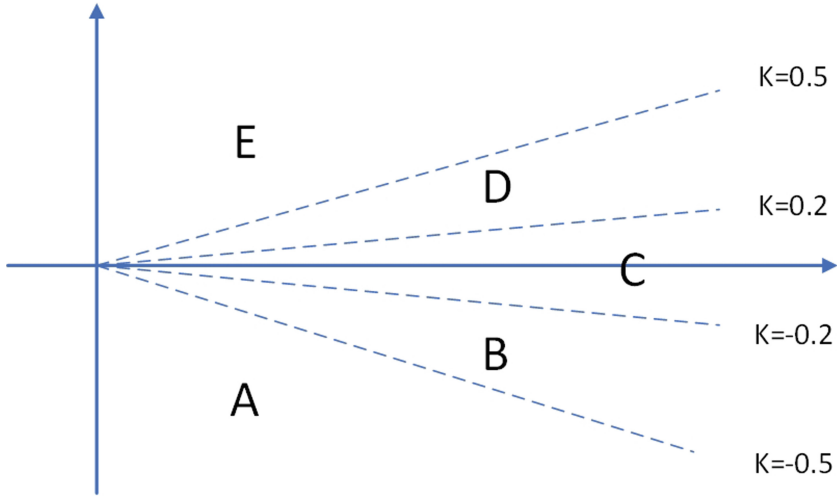


Fig. 3. Schematic diagram of segment slope mapping in the TrSAX method

This method retains the numerical change characteristics and trend information of the original data, which is more conducive to the accurate clustering and grouping of spectrum monitoring data. This method retains the numerical change characteristics and trend information of the original data, which is more conducive to the accurate clustering and grouping of spectrum monitoring data.

3 Spectrum Prediction Based on Seq2seq

Seq2seq model is a typical encoder decoder architecture, which is usually composed of RNN, GRU and LSTM. The encoder encodes the variable length sequence data into a fixed length intermediate vector, and then the decoder decodes the fixed length intermediate vector as the prediction output, realizing that the input sequence of any length can be mapped to the output sequence of any length. Researchers first applied the model to machine translation in the field of NLP. Later, due to its good performance in multi-step prediction, it was also widely used in time series prediction tasks such as transportation [11] and power load [12].

3.1 Long Short-Term Memory Network (LSTM)

Long Short-Term Memory (LSTM) is a widely used temporal recurrent network. LSTM is a gated RNN. The cleverness of LSTM is that the weight of the self-loop is changed by increasing the input threshold, forgetting threshold and output threshold. When the model parameters are fixed, the integration scale at different times can be dynamically changed, thus avoiding the problem of gradient disappearance or gradient explosion when the general RNN is dealing with long-term dependencies (nodes that are far away in the time series). The LSTM cell is shown in Fig. 4.

As shown in Fig. 4, the forget gate determines which information needs to be discarded from the cell. The calculation process is shown in formula (4):

$$\mathbf{f}_t = \sigma(W_f * [\mathbf{h}_{t-1}, x_t] + b_f) \tag{4}$$

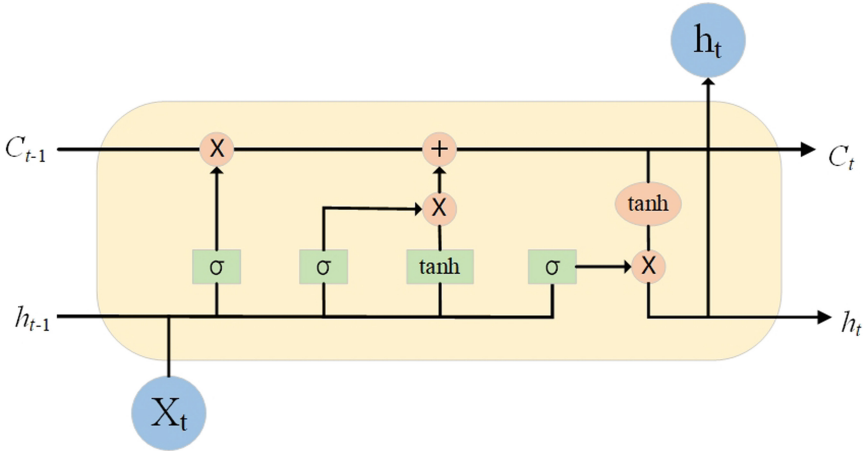


Fig. 4. LSTM cell.

W_f is the weight matrix of the forget gate, where the matrix from cell to gate is diagonal and the rest are non-diagonal. b_f is the offset term of the forget gate and σ is the sigmoid function. The input gate determines what information should be stored. First determine what value needs to be passed through the sigmoid function. The second part is to create a new vector through \tanh function, which will be added to the state C_t . The function of the output gate is to determine the content of the output. The gate processes the cell state through \tanh to get a value between -1 and 1 and multiplies it with the output of the sigmoid layer, and then determines the output part.

Based on the long-term spectrum prediction problem, the LSTM network with long-term memory is obviously very attractive. But at the same time, we also note that the basic LSTM network capability will also decline significantly in the multi output scenario, so the seq2seq architecture has also been introduced into our work.

3.2 Spectrum Prediction Based on Seq2seq

The network architecture of seq2seq based on LSTM used in this paper is shown in Fig. 5. Among them, $X = \{x_1, \dots, x_n\}$ represents the historical monitoring data with step n , and $Y = \{y_1, \dots, y_m\}$ represents the output prediction result with step m . The hidden state h_t at each moment in the encoder is jointly determined by the input data x_t at the current moment, the hidden state h_{t-1} and the cell state c_{t-1} at the previous moment:

$$h_t = f(x_t, h_{t-1}, c_{t-1}) \tag{5}$$

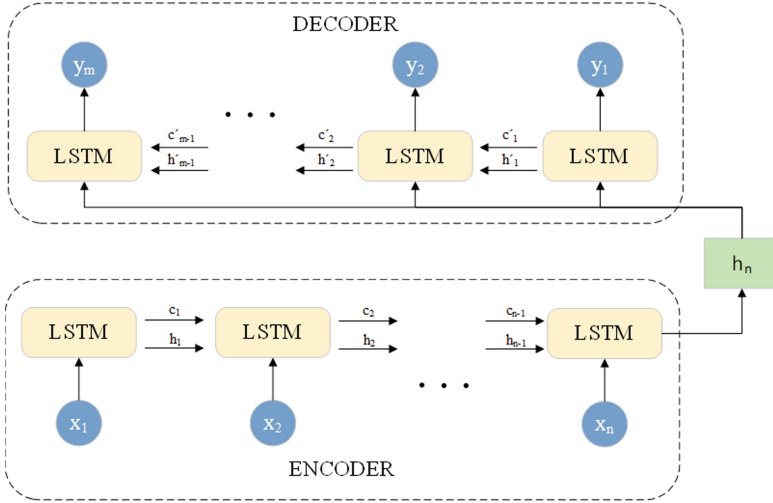


Fig. 5. Network architecture of seq2seq based on LSTM.

The encoder is updated for n time steps, and the input sequence is encoded as the hidden state h_n at the final time step. Due to the long-term memory function of the recurrent neural network, h_n theoretically contains the complete information of the input sequence.

The decoder accepts the final state h_n of the encoder as the initial input value. The hidden state h'_t of the decoder at each moment is updated by the input h_n at the current moment, the hidden state h'_{t-1} and the cell state c'_{t-1} at the previous moment, and its expression is as follows:

$$h'_t = f(h_n, h'_{t-1}, c'_{t-1}) \tag{6}$$

After step-by-step decoding, the final output sequence $Y = \{y_1, \dots, y_m\}$ is formed.

Algorithm: The TrSAX-seq2seq based Spectrum Prediction

Input: Broadband historical monitoring data

Output: Multiple narrowband spectrum prediction models

Begin

1. Normalize the historical observation data of each channel;
 2. Divide the data into N segments of equal length and calculate the mean and slope;
 3. Construct the symbol sequence for each channel using the TrSAX method;
 4. Hierarchical clustering based on symbol sequence to get n channel groupings;
 5. Construct a dataset according to channel grouping and train it to get n seq2seq prediction models
-

End

The pseudocode of the complete flow of this method is given. In short, this paper constructs the feature sequence of the channel based on the trend symbol clustering approximation method and clusters to obtain multiple channel groups, and then trains the seq2seq prediction model separately to realize the long-term spectrum prediction of the wide band.

4 Experimental and Results

In the experimental part, the clustering performance of TrSAX and the final spectrum prediction results are compared and analyzed. In the comparison of clustering, we mainly compare the traditional agglomerative hierarchical clustering methods. In the experiment of spectrum prediction, we compared it with some common methods in the field of spectrum prediction, such as LSTM, CNN and attention mechanism, according to the previous literature. It should be noted that most of these literatures do not disclose their data sets and hyperparameters, so this paper may not perfectly restore their performance, but after enough comparisons, it is enough to illustrate the superiority of our method.

4.1 Datasets and Preprocessing

The data set of the experiment in this paper comes from the open-source spectrum measurement data of Aachen university of technology in Germany [13]. The dataset contains 30 MHz-6 GHz spectrum monitoring data collected from three monitoring sites in Aachen and Maastricht. The frequency resolution is 200 kHz, and the acquisition time interval is about 1.8 s. This paper selects the 1500-3000 MHz monitoring results collected at the third-floor balcony of a residential area in Aachen for about 14 days as the data set of this experiment.

In this paper, the spectrum occupancy sequence is predicted, and the monitoring data is processed into an occupancy sequence with an interval of 15 min according to the technical specifications recommended by ITU. At the same time, in order to avoid the influence of a large number of long-term unused or long-term used channels, the above-mentioned channels are excluded from the experiments in this paper. Because the prediction of these channels is easy and meaningless, it will affect the evaluation of algorithm performance. In the prediction experiment in this paper, we use the historical data of the past two days to predict the outcome of the next day. That is, the input sequence length of our prediction model is 192, and the output result length is 96.

4.2 Analysis of Results

We first explore the impact of the TrSAX method on the clustering effect of spectrum monitoring data. The clustering evaluation index used in this paper is the silhouette coefficient. This index is an evaluation method of clustering effect, which combines two factors: cohesion and separation. It can be used to evaluate the impact of different algorithms or different operating modes of algorithms on the clustering results on the

basis of the same original data. The results are shown in Fig. 6. We mainly focus on the comparison of the highest points between the curves, that is, the best performance of the two clustering methods under the optimal number of clusters. It can be seen that the best performance of the TrSAX method occurs when the number of clusters is 5, and its silhouette coefficient is much higher than the result of directly clustering the original sequence under all the number of clusters. In contrast, the best performance for clustering using the original sequence occurs when the number of clusters is 2. Obviously, the result of clustering into 2 clusters does not group the data very well, which shows that the clustering results using the original sequence are far from ideal. In the subsequent comparative experiments, the number of clusters used in this paper is uniformly 5.

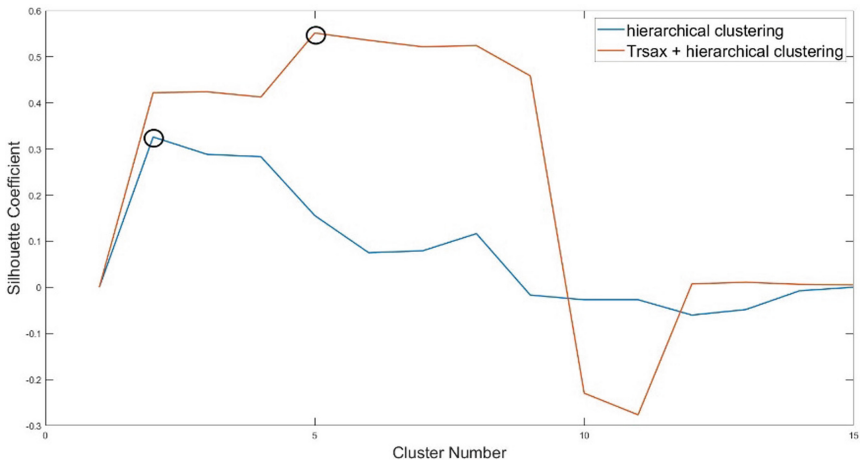


Fig. 6. Performance based on TrSAX clustering.

Then, we compare the performance of spectrum prediction. In this experiment, we uniformly use three layers of LSTM, and the number of hidden units in each layer of LSTM is 128. Experiments were carried out on five clusters based on TrSAX and hierarchical clustering. The experimental results are shown in Fig. 7. It can be seen that the RMSE of multiple networks in the first two clusters are almost the same, and the clustering based method in the last three clusters has achieved obvious advantages, of which our method has more obvious advantages. It can be predicted that this advantage will become more obvious as the predicted spectrum bandwidth increases. Therefore, we believe that clustering grouping is a simple and effective method to improve the performance of broadband prediction involving multiple services.

Finally, we compared the performance of different methods on all prediction steps, as shown in Fig. 8. As shown in the blue curve in the figure, the trsax-seq2seq method used in this paper achieves the overall minimum root mean square error, and the error is relatively stable in all steps. The performance of other clustering based methods is second. The performance of non clustering method is average. We believe that, in the case of the broadband data set used in this paper, even if the clustering method is not

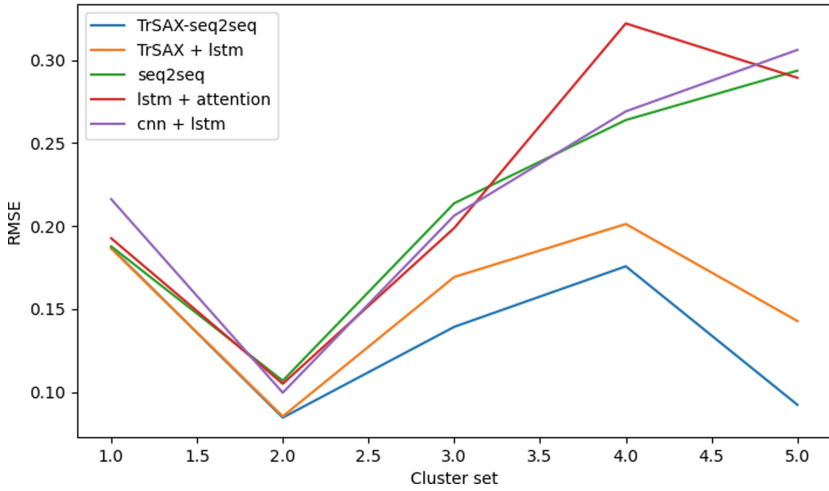


Fig. 7. RMSE performance comparison of different models in five groups

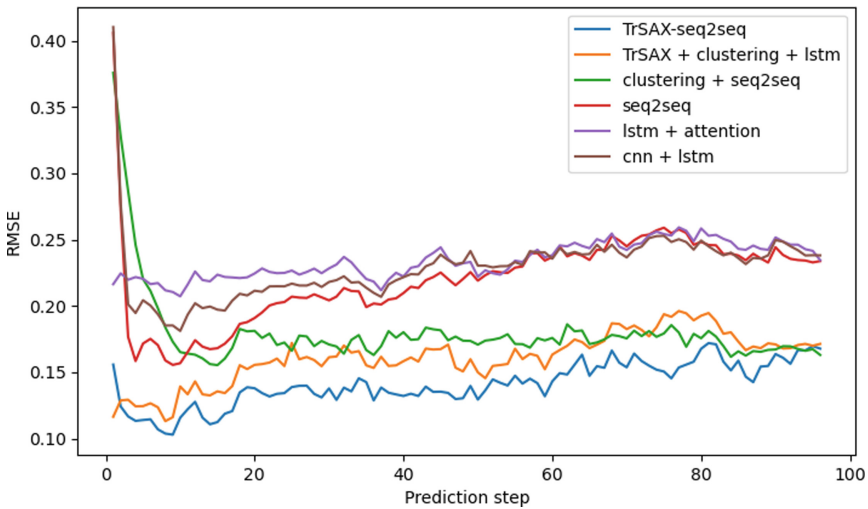


Fig. 8. Performance comparison of different methods on all prediction steps.

ideal, its performance index has been significantly improved. It further shows that our method is effective in this broadband prediction scenario.

5 Conclusion

In this paper, the problem of long-term spectral prediction over broadband is investigated. Based on the difficulties faced by this problem, a clustering and grouping method of spectral data based on TrSAX expression was developed and combined with the seq2seq

network to improve the accuracy of prediction. Experiments are carried out on a real open source spectrum monitoring dataset, and the experimental results demonstrate the effectiveness of the method proposed in this paper. We expect to introduce a more advanced encoder-decoder structure model Transformer applied to this problem in the future to further improve the prediction performance.

Acknowledgment. This work is supported by the National Natural Science Foundation of China (61771154) and the Fundamental Research Funds for the Central Universities (3072021CF0801).

This work is also supported by Key Laboratory of Advanced Marine Communication and Information Technology, Ministry of Industry and Information Technology, Harbin Engineering University, Harbin, China.

References

1. Lin, Y., Wang, C., Wang, J., Dou, Z.: A novel dynamic spectrum access framework based on reinforcement learning for cognitive radio sensor networks. *Sensors* **16**(10), 1675 (2016)
2. Guo, L., Wang, M., Lin, Y.: Electromagnetic environment portrait based on big data mining. *Wireless Communications and Mobile Computing* (2021)
3. Mosavat-Jahromi, H., Li, Y., Cai, L., Pan, J.: Prediction and modeling of spectrum occupancy for dynamic spectrum access systems. *IEEE Transactions on Cognitive Communications and Networking* **7**(3), 715–728 (2021)
4. Wang, Y., Zhang, Z., Ma, L., Chen, J.: SVM-based spectrum mobility prediction scheme in mobile cognitive radio networks. *The Scientific World Journal* (2014)
5. Wang, X., Peng, T., Zuo, P., Wang, X.: Spectrum Prediction Method for ISM Bands Based on LSTM. In: 2020 5th International Conference on Computer and Communication Systems (ICCCS), pp. 580–584. IEEE (2020 May)
6. Shawel, B.S., Woldegebreal, D.H., Pollin, S.: Convolutional LSTM-based long-term spectrum prediction for dynamic spectrum access. In: 2019 27th European Signal Processing Conference (EUSIPCO), pp. 1–5. IEEE (2019 September)
7. Shi, C., Dou, Z., Lin, Y., Li, W.: Dynamic threshold-setting for RF-powered cognitive radio networks in non-Gaussian noise. *Physical Communication* **27**, 99–105 (2018)
8. Yu, Y., Zhu, Y., Wan, D., Liu, H., Zhao, Q.: A novel symbolic aggregate approximation for time series. In: Lee, S., Ismail, R., Choo, H. (eds.) *IMCOM 2019. AISC*, vol. 935, pp. 805–822. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-19063-7_65
9. Yahyaoui, H., Al-Daihani, R.: A novel trend based SAX reduction technique for time series. *Expert Syst. Appl.* **130**, 113–123 (2019)
10. Sun, J., Wang, J., Chen, J., Ding, G., Lin, F.: Clustering analysis for internet of spectrum devices: real-world data analytics and applications. *IEEE Internet Things J.* **7**(5), 4485–4496 (2020)
11. Zhang, Z., Li, M., Lin, X., Wang, Y., He, F.: Multistep speed prediction on traffic networks: A deep learning approach considering spatio-temporal dependencies. *Transportation research part C: emerging technologies* **105**, 297–322 (2019)
12. Gong, G., An, X., Mahato, N.K., Sun, S., Chen, S., Wen, Y.: Research on short-term load prediction based on Seq2seq model. *Energies* **12**(16), 3199 (2019)
13. Wellens, M., Mähönen, P.: Lessons learned from an extensive spectrum occupancy measurement campaign and a stochastic duty cycle model. *Mobile networks and applications* **15**(3), 461–474 (2010)
14. Zhang, L., Jia, M.: Accurate Spectrum Prediction Based on Joint LSTM with CNN toward Spectrum Sharing. In: 2021 IEEE Global Communications Conference (GLOBECOM), pp. 1–6. IEEE (2021 December)