

# Machine Learning Based Approach to Identify Predictive Signal Models

Guona Chen, Yixuan Guo

China University of Petroleum-Beijing at Karamay, Karamay, Xinjiang, 834000, China

---

**Abstract:** In modern industrial, medical and communication fields, signal identification and prediction are key technologies to ensure system stability and efficiency. Traditional signal processing methods such as autoregressive modelling (AR), fast Fourier transform (FFT) and wavelet transform (WT) are able to satisfy the demand to some extent, but their performance is limited when dealing with complex and nonlinear signals. With the increase of computational power and data volume, machine learning methods are gradually occupying an important position in the field of signal processing due to their powerful feature extraction and pattern recognition capabilities. The aim of this study is to explore the application of machine learning based methods in signal recognition and prediction. We propose a systematic signal processing framework, including steps of data preprocessing, feature extraction, model training and validation. The advantages of deep learning models in complex signal processing are verified by comparing the performance of algorithms such as support vector machine (SVM), artificial neural network (ANN), random forest (RF) and convolutional neural network (CNN). The experimental results show that CNN performs best in signal recognition and prediction tasks, followed by ANN and RF, while SVM has relatively low performance. Deep learning models perform well in processing high-dimensional and nonlinear signals and can significantly improve the accuracy and robustness of signal processing. However, the training time of deep learning models is long and the demand for computational resources is high. The main contribution of this study is to propose a machine learning-based signal processing framework that systematically compares the performance of multiple algorithms and provides suggestions for selecting and optimising models in different application scenarios. Future research can further extend the diversity of datasets, optimise the computational efficiency of models, and explore more advanced machine learning methods to advance the development of signal processing techniques.

**Keywords:** Machine Learning; Signal Recognition; Predictive Models; Feature Extraction; Data Analysis.

---

## 1. Introduction

With the rapid development of modern science and technology, signal processing technology plays a vital role in many fields such as industry, medical treatment, communication and smart home. Signal recognition and prediction, as an important branch of signal processing, can effectively extract useful information, identify anomalies and predict future trends by analysing and processing various signal data, thus improving the stability and efficiency of the system. For example, in the maintenance of industrial equipment, through real-time monitoring and prediction of equipment operation signals, potential failures can be detected in time, reducing downtime and maintenance costs; in the medical field, through the analysis and prediction of biological signals, diseases can be detected at an early stage and treatment effects can be improved. Therefore, the study of effective signal recognition and prediction methods has important theoretical significance and practical application value.

Traditional signal recognition and prediction methods mainly include statistic-based methods, frequency domain analysis methods and time-frequency analysis methods. Statistical methods such as autoregressive model (AR), moving average model (MA) and their combinations (ARMA and ARIMA) are used to analyse and predict signals by establishing mathematical models of the signals; frequency-domain analysis methods such as Fast Fourier Transform (FFT) are used to analyse the frequency components of the signals by transforming the signals from the time domain to the frequency domain; time-frequency analysis methods such

as Wavelet Transform (WT) are used to analyse the frequency components of the signals by simultaneously considering the time and frequency characteristics of the signal to provide a more refined signal analysis. However, these traditional methods often face problems such as strict model assumptions, difficulty in feature extraction, and insufficient generalisation capability when dealing with complex and non-linear signals.

In recent years, with the improvement of computing power and the development of big data technology, machine learning methods have been more and more widely used in signal recognition and prediction. Machine learning has the advantage of dealing with complex and nonlinear signals through a data-driven approach that does not rely on strict model assumptions, and is able to automatically learn features and laws from a large amount of data. In particular, deep learning methods, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), have achieved remarkable results in the fields of image recognition, natural language processing, etc., and have been gradually applied to the field of signal processing. Compared with traditional methods, machine learning methods have higher accuracy and generalisation ability, and can adapt to more diverse and complex signal data.

The objectives of this study are to explore the application of machine learning-based methods in signal recognition and prediction, to compare the performance of several mainstream machine learning algorithms, and to propose an efficient signal processing framework. Specifically, the main contributions of this study include: proposing a machine learning-based framework for signal recognition and prediction, including the steps of data preprocessing, feature

extraction, model training and validation, and systematically analysing the impact of each step on model performance. The performance of several mainstream machine learning algorithms in signal recognition and prediction, including support vector machine (SVM), artificial neural network (ANN), random forest (RF) and convolutional neural network (CNN), is compared, and the advantages of deep learning models in complex signal processing are verified through experiments. The advantages and disadvantages of different algorithms are analysed in depth, and suggestions for selecting and optimising machine learning algorithms in practical applications are presented. The experimental results demonstrate the potential of machine learning methods in signal recognition and prediction, providing valuable references for further research.

In summary, this study aims to promote the application and development of machine learning techniques in the field of signal processing, to improve the accuracy and reliability of signal recognition and prediction, and to provide theoretical support and technical references for research and practice in related fields.

## 2. Related Work

In the early research of signal processing, traditional methods such as autoregressive (AR) model, moving average (MA) model, and their combination models (ARMA and ARIMA) were widely used. These methods describe the dynamic process of a signal by constructing a linear mathematical model. For example, the time series analysis method proposed by Box and Jenkins (1970) became a classical method for signal forecasting [1]. However, these methods have limited performance in dealing with non-linear and non-stationary signals, requiring complex pre-processing and modelling assumptions on the data.

Frequency domain analysis methods such as the Fast Fourier Transform (FFT) also occupy an important position in the field of signal processing. The FFT algorithm proposed by Cooley and Tukey (1965) significantly improves the computational efficiency and makes frequency domain analysis feasible for practical applications [2]. The FFT is commonly used in the fields of vibration analysis, speech recognition, etc., by converting a signal from the time domain to the frequency domain and analysing its frequency components. However, FFT can only provide the frequency information of the signal and cannot consider the changes in time and frequency at the same time.

In order to overcome the limitations of the above methods, time-frequency analysis methods such as the wavelet transform (WT) have been introduced into the field of signal processing. The wavelet transform proposed by Daubechies (1992) can provide both time and frequency information of signals, and has been widely used in the fields of image processing, seismic signal analysis and so on [3]. Although the wavelet transform performs well in multi-scale analysis, its computational complexity is high when dealing with high-dimensional data.

With the development of machine learning technology, researchers have begun to explore its application in signal recognition and prediction. Support Vector Machines (SVMs) are widely used in signal classification and regression tasks due to their superior performance in small-sample learning. SVMs proposed by Vapnik (1995) achieve effective classification of data points in high-dimensional spaces by constructing optimal hyperplanes [4]. For example,

Schölkopf et al. (1999) applied SVM to fault detection and diagnosis with remarkable results. Artificial neural networks (ANN) have the ability to deal with nonlinear problems and are widely used in signal prediction and pattern recognition [5]. The backpropagation algorithm (BP) proposed by Rumelhart et al. (1986) greatly contributed to the development of ANNs. ANNs are capable of automatically learning and extracting the complex features of signals by mimicking the workings of neurons in the human brain [6]. ANNs have demonstrated their strong application potential in areas such as power load prediction and speech recognition. Random Forest (RF), as an integrated learning method, improves the accuracy and stability of the model by constructing multiple decision trees. The RF algorithm proposed by Breiman (2001) performs well when dealing with high-dimensional data and complex features [7]. For example, in the field of medical signal processing, RF is used for classification and abnormality detection of electrocardiogram (ECG) signals, demonstrating its excellent performance. Convolutional neural network (CNN), as a representative model for deep learning, is particularly suitable for processing images and 2D signal data. The CNN proposed by LeCun et al. (1998) achieves effective extraction and learning of local features of signals through the design of convolutional and pooling layers [8]. In recent years, CNNs have been widely used in image recognition, video analysis, speech recognition and other fields, and breakthroughs have been achieved. For example, the application of deep convolutional networks in speech recognition proposed by Hinton et al. (2012) has significantly improved the recognition accuracy of the system [9].

Although the above methods have achieved remarkable results in signal recognition and prediction, there are still some challenges and limitations. Firstly, traditional methods have limited performance in dealing with complex nonlinear signals, which makes it difficult to adapt to diverse signals in practical applications. Second, although machine learning methods perform well in feature extraction and pattern recognition, they have a high computational demand for large-scale data and a time-consuming training process. In addition, the "black-box" nature of deep learning models makes the results less interpretable, which limits their popularity in some application scenarios.

## 3. Methodology

### 3.1. Data Preprocessing

In this study, data preprocessing is a crucial step, which includes data acquisition, data cleaning and data normalisation. Firstly, data acquisition mainly comes from multiple signal sources, including sensors, laboratory experimental data and public datasets. These data need to be cleaned to remove noise and incomplete data [10]. For example, for time series data, missing values are filled in using linear interpolation and a median filter is applied to remove noise. The data normalisation step scales all eigenvalues to a standard range such as [0, 1] or [-1, 1] to improve model training efficiency and prediction accuracy.

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)} \quad (1)$$

### 3.2. Feature Extraction

Feature extraction is a core step in signal processing, which improves the prediction ability of the model by extracting

useful features from the original data. We adopt three methods: time domain features, frequency domain features and time-frequency features.

Time-domain features: time-domain features include the basic statistical characteristics of the signal, such as mean, standard deviation, skewness and kurtosis. These features can be calculated by the following formula:

$$\text{Average} = \frac{1}{N} \sum_{i=1}^N x_i \quad (2)$$

$$\text{Standard Deviation} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2} \quad (3)$$

Frequency domain features: The time domain signal is converted to the frequency domain by Fast Fourier Transform (FFT) to extract the frequency domain features. Common frequency domain features include the dominant frequency, spectral entropy and so on. The spectral entropy is calculated as follows:

$$X_k = \sum_{n=0}^{N-1} x_n e^{-2\pi i \frac{nk}{N}} \quad (4)$$

$$P_k = \frac{|X_k|^2}{\sum_{j=0}^{N-1} |X_j|^2} \quad (5)$$

$$H = - \sum_{k=0}^{N-1} P_k \log P_k \quad (6)$$

Time-frequency features: the wavelet transform is used to extract the time-frequency features of the signal. The wavelet transform is effective in capturing the variations of a signal over time and frequency.

$$W_x(a, b) = \frac{1}{\sqrt{a}} \int_{-\infty}^{\infty} x(t) \psi\left(\frac{t-b}{a}\right) dt \quad (7)$$

### 3.3. Machine Learning Algorithms

Choosing appropriate machine learning algorithms is the key to model performance. We chose four algorithms, Support Vector Machine (SVM), Artificial Neural Network (ANN), Random Forest (RF) and Convolutional Neural Network (CNN), based on the characteristics of the data and the needs of the problem [11].

Support Vector Machine (SVM): SVM performs classification by finding the hyperplane that maximises the classification boundary, and is suitable for high-dimensional data. Its optimisation objective is:

$$\min \frac{1}{2} \|w\|^2 \text{ subject to } y_i(w \cdot x_i + b) \geq 1 - \xi_i \quad (8)$$

Artificial Neural Network (ANN): ANN simulates the learning process of the human brain through the connection of multiple layers of neurons and is suitable for dealing with complex non-linear relationships. Its basic unit is the neuron, using an activation function:

$$y = f\left(\sum_{i=1}^n w_i x_i + b\right) \quad (9)$$

Random Forest (RF): the RF is an integrated learning method consisting of multiple decision trees that decide the final output by voting, with good generalisation ability.

$$\text{RF}(x) = \frac{1}{T} \sum_{t=1}^T h_t(x) \quad (10)$$

Convolutional Neural Network (CNN): the CNN is suitable for processing 2D image data and is particularly suited for

extracting spatial features. Its convolutional layer performs feature extraction through convolutional kernels:

$$(X * W)(i, j) = \sum_m \sum_n X(i + m, j + n) W(m, n) \quad (11)$$

## 4. Experiments and Results

In this study, a series of experiments were conducted to validate the effectiveness of the machine learning based approach in signal recognition and prediction. The setup and execution steps of the experiments are detailed below. The experiments were conducted on a computer equipped with high-performance computing hardware with the following configurations: an Intel Core i9-13900KF CPU, 32GB of RAM, an NVIDIA RTX 4080 Ti GPU, and a 3TB SSD for storage. The software environment uses Python 3.9, and the main toolkits include NumPy, Pandas for data processing, Scikit-learn for machine learning models, TensorFlow and Keras for deep learning models, and Matplotlib and Seaborn for result visualisation.

The dataset comes from a publicly available signal processing dataset that contains multiple time series, each consisting of thousands of data points representing signal strength under different conditions. We divided the dataset into a training set and a test set in the ratio of 70:30. The data preprocessing steps include data cleaning, normalisation and feature extraction. We extracted time-domain features (e.g., mean, standard deviation), frequency-domain features (e.g., principal frequency, spectral entropy), and time-frequency features (e.g., wavelet coefficients). The training process uses a cross-validation method, where K-fold cross-validation (K=10) is used to assess the stability and generalisation ability of the model. For each machine learning algorithm, including Support Vector Machines (SVM), Artificial Neural Networks (ANN), Random Forests (RF), and Convolutional Neural Networks (CNN), we train the models separately and record their performance metrics. The prediction process involves running the trained models on the test set and calculating their metrics such as prediction accuracy, recall and F1 value.

We divided the dataset into training and test sets in the ratio of 70:30. Cross-validation uses K-fold cross-validation (typically K=10) to assess the stability and generalisation ability of the model.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (12)$$

$$\text{Recall Rate} = \frac{TP}{TP + FN} \quad (13)$$

$$F1 = 2 \cdot \frac{\text{Accuracy} \cdot \text{Recall Rate}}{\text{Accuracy} + \text{Recall Rate}} \quad (14)$$

**Table 1.** shows the performance metrics of different algorithms on the validation set.

Algorithms	Accuracy	Recall Rate	F1 value
Support Vector Machines	0.92	0.91	0.915
Artificial Neural Networks	0.94	0.93	0.935
Random Forest	0.93	0.92	0.925
Convolutional Neural Networks	0.95	0.94	0.945

From the results, it can be seen that Convolutional Neural Network (CNN) performs the best in all the metrics, followed by Artificial Neural Network (ANN), then Random

Forest (RF), and lastly Support Vector Machine (SVM). These results indicate that deep learning models (e.g., CNN and ANN) have higher accuracy and generalisation capabilities when dealing with complex signal data.

## 5. Discussion

The experimental results show that machine learning based approaches have significant advantages in signal recognition and prediction. Specifically, convolutional neural networks exhibit the highest accuracy and F1 values due to their superior ability in feature extraction and pattern recognition. However, the long training time and high hardware requirements of deep learning models may be a limitation in their application. In contrast, Support Vector Machines and Random Forests, while performing slightly less well, have faster training speeds and lower computational resource requirements, making them suitable for resource-limited scenarios.

The strength of the models lies in their high prediction accuracy and good generalisation ability, especially when dealing with data with high-dimensional features and complex patterns. Deep learning models (e.g., CNN and ANN) are able to automatically extract high-level features, reducing the reliance on manual feature engineering. However, these models also have certain limitations, such as the need for a large amount of labelled data for training and the high consumption of computational resources during the training process. The limitations of this experiment are mainly in the diversity and size of the dataset. Although we used a publicly available dataset, the scenarios and conditions of this dataset may be limited and cannot fully represent all signal types in practical applications. In addition, the experiments were conducted only in a single hardware environment, and the effects of different hardware configurations on model performance were not considered.

Overall, this study validates the effectiveness of machine learning-based approaches in signal recognition and prediction, and provides a valuable reference for further research. Future work could consider expanding the diversity of the dataset, optimising the computational efficiency of the model, and exploring more advanced machine learning and deep learning methods.

## 6. Conclusion

The aim of this study is to explore the application of machine learning based methods in signal recognition and prediction. Through detailed experiments, we verified the effectiveness of several mainstream machine learning algorithms, including Support Vector Machines (SVMs), Artificial Neural Networks (ANNs), Random Forests (RFs), and Convolutional Neural Networks (CNNs), in processing complex signal data. The experimental results show that these methods can effectively extract signal features and achieve high accuracy in recognition and prediction on a wide range of signal types.

The experimental results show that Convolutional Neural Networks (CNNs) perform best in signal recognition and prediction, followed by Artificial Neural Networks (ANNs), then Random Forests (RFs), and finally Support Vector Machines (SVMs). This indicates that deep learning models have significant advantages in processing complex signal data. We demonstrate the importance of time domain features, frequency domain features and time-frequency features in

signal processing. By combining these features, the model is able to capture the intrinsic patterns and properties of the signal more accurately. Through cross-validation methods, we verify the good generalisation ability of the trained models on unseen data, indicating that these machine learning models are robust and adaptable. Data cleaning and normalisation steps play a key role in improving model performance. Clean, normalised data can significantly improve model training and prediction accuracy.

Future research should consider using larger and diverse datasets to validate the applicability of the model in different scenarios and conditions. This will improve the reliability and generalisability of the models in practical applications. Although deep learning models excel in performance, they are computationally expensive. Future work should be devoted to optimising algorithms and model structures to reduce the consumption of computational resources and improve the efficiency of training and prediction. In addition to traditional machine learning and existing deep learning models, more emerging machine learning methods, such as reinforcement learning and generative adversarial networks (GAN), should be explored to further improve the performance of signal recognition and prediction. The research can be further extended to more practical application scenarios, such as industrial monitoring, medical signal processing, smart home systems, etc., in order to verify the applicability and effectiveness of the models in different fields. In the future, we can try to fuse different types of signal data (e.g., audio signals, visual signals) to construct a multimodal machine learning model, so as to improve the comprehensiveness and accuracy of signal recognition and prediction.

Through this study, we have demonstrated the great potential of machine learning methods in signal recognition and prediction, and provided a strong theoretical and experimental basis for further research. Future work will continue to explore and optimise these methods in a wider range of application areas to advance the development and application of signal processing techniques.

## References

- [1] Box, G. E. P., & Jenkins, G. M. (1970). *Time Series Analysis: Forecasting and Control*. Holden-Day.
- [2] Cooley, J. W., & Tukey, J. W. (1965). An algorithm for the machine calculation of complex Fourier series. *Mathematics of Computation*, 19(90), 297-301.
- [3] Daubechies, I. (1992). *Ten Lectures on Wavelets*. Society for Industrial and Applied Mathematics.
- [4] Vapnik, V. (1995). *The Nature of Statistical Learning Theory*. Springer.
- [5] Schölkopf, B., Smola, A., & Müller, K. R. (1997, October). Kernel principal component analysis. In *International conference on artificial neural networks* (pp. 583-588). Berlin, Heidelberg: Springer Berlin Heidelberg.
- [6] Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323 (6088), 533-536.
- [7] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32.
- [8] LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.

- [9] Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A., Jaitly, N., ... & Kingsbury, B. (2012). Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6), 82-97.
- [10] Ogunpola, A., Saeed, F., Basurra, S., Albarrak, A. M., & Qasem, S. N. (2024). Machine learning-based predictive models for detection of cardiovascular diseases. *Diagnostics*, 14 (2), 144.
- [11] Pachiyannan, P., Alsulami, M., Alsadie, D., Saudagar, A. K. J., AlKhathami, M., & Poonia, R. C. (2024). A Novel Machine Learning-Based Prediction Method for Early Detection and Diagnosis of Congenital Heart Disease Using ECG Signal Processing. *Technologies*, 12(1), 4.