

## A Novel Power Load Identification Strategy Based on CEEMD-PCA-AE-Shapelet Dictionary Learning Algorithm

Junxiong Ge<sup>1\*</sup>, Mingmin Yuan<sup>2</sup>, Zhu Tang<sup>2</sup>, Quancheng Pan<sup>2</sup>, Congsu Jin<sup>2</sup>, Haimin Hong<sup>1</sup>

<sup>1</sup>Shenzhen China Gridcom Technology Communication Co., Ltd, 518028, Shenzhen, China

<sup>2</sup>State Grid Beijing Electric Power Company, Beijing, 100031, China

### Abstract

The development of the Energy Internet under the China 'dual-carbon' strategy has brought new challenges in load identification and anomaly detection, particularly in industrial applications. Traditional methods often struggle with the sparse, noisy nature of industrial load data, limiting their generalization ability and accuracy. This paper proposes a novel load identification method based on the CEEMD-PCA-AE-Shapelet dictionary learning algorithm, which addresses these challenges by combining an improved unsupervised anomaly detection model (CEEMD-PCA-AE) with the Shapelet dictionary learning algorithm. Firstly, CEEMD-PCA-AE enhances the model's ability to handle non-linear and noisy data, significantly improving the generalization and accuracy of industrial load anomaly detection. Secondly, the Shapelet dictionary learning algorithm reduces computational complexity and improves model performance by incorporating dictionary learning into time-series classification. The proposed method outperforms traditional models, as demonstrated by numerical experiments, offering enhanced load identification accuracy, and improved detection of anomalies. This method has significant potential for optimizing the efficiency of energy management in industrial settings and can be applied globally, as shown by its effectiveness on real-world data from Liaoning Province, China.

**Keywords:** power data, energy internet, load identification, data detection, learning algorithm

Received on 5 November 2025, accepted on 14 December 2025, published on 31 March 2026

Copyright © 2026 Junxiong Ge *et al.*, licensed to EAI. This is an open access article distributed under the terms of the [CC BY-NC-SA 4.0](#), which permits copying, redistributing, remixing, transformation, and building upon the material in any medium so long as the original work is properly cited.

doi: 10.4108/ew.11913

### 1. Introduction

Under the China “dual-carbon” strategy, power system transformation has attracted much attention [1-2], and the grid integration of renewable energy and controllable loads has brought new challenges to power system management, especially in load identification and anomaly detection [3-5]. Industrial power consumption accounts for a large portion of the total power consumption, its time pattern is complex [6], and the data is usually characterized by high dimensionality, sparsity, and noise, etc [7-8]. Traditional load identification

methods are difficult to deal with these characteristics, which leads to a reduction in load identification efficiency [9-10].

Anomaly detection and unknown load identification play a crucial role in understanding industrial electricity consumption [11]. However, detecting anomalies becomes increasingly difficult as the volume and complexity of data grows. Currently, deep learning-based models are widely used in this field of anomaly detection [12-13]. Deep learning models such as Deep Belief Networks, Convolutional Neural Networks and Autoencoders have demonstrated strong capabilities in processing high-dimensional features and detecting anomalies [14-15]. However, due to the presence of

\*Corresponding author. Email: 2125786732@qq.com

sparsity and noise in industrial load data, traditional methods still face challenges and need to be improved to increase their efficiency [16]. In addition, traditional load identification techniques usually rely on direct feature extraction, processing each sampling point independently while ignoring the time dependencies. Since industrial load profiles exhibit strong time-series characteristics, disregarding these dependencies can lead to poor classification accuracy and interpretability [17]. In addition, sampling processes in industrial environments occur at intervals of 15-30 minutes, sometimes more frequently, which can lead to increased data redundancy and computational complexity. Therefore, there is a need to propose improvements to the traditional load identification methods in order to improve the efficiency and quality of load identification.

To address these challenges, a new energy Internet load identification method based on CEEMD-PCA-AE-Shapelet dictionary learning algorithm is proposed in this paper. The CEEMD-PCA-AE in the load identification method proposed in this paper solves the problem of noise and outliers, and the Shapelet dictionary learning algorithm is able to classify and identify unknown loads. Firstly, CEEMD is combined with PCA and AE for the problems of poor model generalization and the inability to deal with non-linear complex data. Secondly, the idea of dictionary is introduced into the Shapelet algorithm in order to reduce the computation and improve the generalization ability of the model. Finally, a novel energy internet load identification method based on CEEMD-PCA-AE-Shapelet dictionary learning algorithm is proposed to optimize the detection of abnormal load data and to improve the computing speed and generalization ability of load identification.

The contributions of this paper are summarized as follows:

(1) A CEEMD-PCA-AE-Shapelet dictionary learning model is proposed to improve the accuracy of load identification for the first time. The model is used to detect anomalies in industrial load data and optimize unknown load identification. Thus, the accuracy of load recognition can be improved. This results in better and accurate identification of unknown load types. The experimental results demonstrate that the CEEMD-PCA-AE-Shapelet dictionary learning model improves the accuracy of load identification.

(2) The method saves the characteristics of target modelling data and simulation requirements. The CEEMD-PCA-AE-Shapelet dictionary learning model is built to detect abnormal load data. The optimization method based on the introduction of the dictionary learning idea improves the generalization ability of the load recognition model. This enables the reconstructed load recognition model to effectively identify the type of unknown load.

(3) This paper fills the gap in the literature by investigating the trend of the impact of different load identification methods on the performance of industrial load identification. The error of including large fluctuation data in the overall load model modelling is reduced. The combination of CEEMD-PCA-AE and Shapelet dictionary learning algorithm effectively improves the load identification accuracy, and a CEEMD-PCA-AE-Shapelet dictionary learning model is proposed. The

effectiveness of the algorithm is verified by the actual load data in Liaoning Province, China.

The rest of the paper is organized as follows. Section 2 introduces the CEEMD-PCA-AE model and discusses how it can solve the problem of sensitivity to noise and outliers. Section 3 describes the Shapelet dictionary learning algorithm and how it improves the running speed and generalization of load recognition. A case study is presented in Section 4. Finally, Section 5 gives conclusions.

## 2. Rationale and role of the CEEMD-PCA-AE model

CEEMD is a method used for the decomposition of nonlinear and non-smooth signals, which decomposes the original signal into multiple intrinsic mode functions (IMFs), which in turn allows the extraction of key feature information from the signal [18-19]. PCA is a method for dimensionality reduction of data, which aims at identifying the main components of the data for the purpose of dimensionality reduction [20]. AE is a neural network model that learns the latent representations in the data and reconstructs the original data [21]. Combining CEEMD, PCA and AE can effectively detect anomalies in industrial load data.

Firstly, since the traditional PCA method is very sensitive to noise and often interfered by noise, which leads to the output results not matching with the actual situation, the CEEMD method is used to reduce the noise of the data set in order to improve the model's anti-noise performance. Secondly, the self-encoder can achieve excellent outlier detection effect when dealing with low-dimensional data, however, when faced with sparse data for dimensionality reduction, its detection performance will be significantly reduced, and it is difficult to achieve effective identification of outliers. In view of the shortcomings of the current data anomaly detection model based on the conventional autoencoder, it is improved by introducing the idea of principal component analysis (PCA) in the construction of the model, and dividing the input data into normal and abnormal data parts, in which the normal data part will be reconstructed and outputted by the autoencoder. The PCA is more effective in dealing with the linear data, but the detection effect is average when facing the nonlinear data, and the AE has a strong capability of dealing with the nonlinear data, so the two are combined to improve the accuracy of anomaly detection.

The CEEMD-PCA-AE model is divided into four main stages:

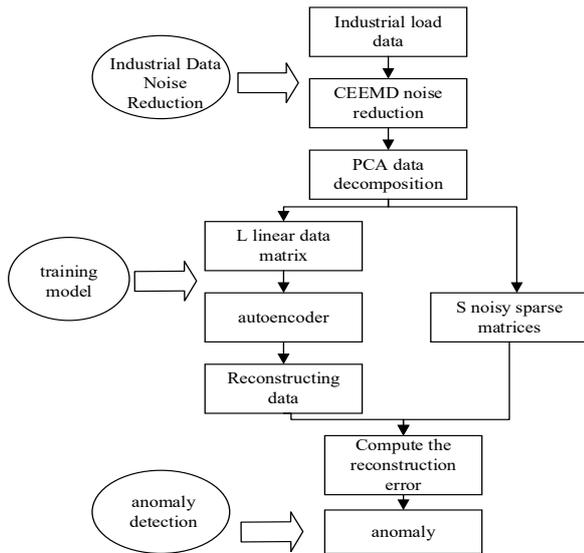


Figure 1. CEEMD-PCA-AE mode

(1) CEEMD performs noise reduction. (2) PCA is responsible for selecting key features from the data to reduce the data dimensionality and removing redundant features, and then this dimensionality reduced data is used as an input for AE; (3) the input data is reconstructed by using the AE method, and then it is coded and decoded to obtain the reconstructed data; and (4) the reconstruction error is calculated and outliers are judged. Figure 1 shows the CEEMD-PCA-AE anomaly detection model. Figure 1 shows the CEEMD-PCA-AE outlier detection model.

### 3. Industrial Load Recognition Based on Shapelet Dictionary Learning Algorithm

#### 3.1 Shapelet Algorithm

Shapelet refers to the most recognizable sequence in a time series, which can maximize the category information of the time series. The distance between each sampling point and Shapelet is calculated as a criterion of similarity, which is used as a load recognition time series feature. On this basis, a new load identification method is proposed, i.e., whether the time series samples contain Shapelets of corresponding categories is used to achieve unknown load identification. The time-series trajectory features fully consider the continuity of the sample points, which can effectively improve the classification accuracy. The use of Shapelet method to extract time-series trajectory features with low dimensionality can effectively improve the classification efficiency. The working principle of Shapelet is as follows:

(1) Traversing the data and extracting candidate Shapelet. each time series in the loaded dataset is traversed using sliding window and time series with similar characteristics are extracted as candidate Shapelet.

(2) Calculate Information Gain (IG). Information gain is constructed by calculating the shortest distance between the extracted Shapelet and each time series to construct the shortest distance ranking line. To calculate the Information Gain, the maximum split point of IG is selected as the optimal split point, as shown in equation (1).

$$IG = E(T) - \frac{n_1}{n} E(T_{left}) - \frac{n_2}{n} E(T_{right}) \quad (1)$$

where  $n_1$  and  $n_2$  are the number of time series of category 1 and category 2, respectively,  $T_{left}$  and  $T_{right}$  are the sets composed of time series to the left and right of the split point, respectively, and  $E(\cdot)$  is the information entropy function, as shown in equation (2).

$$E(D) = -\sum_{i=1}^2 \frac{n_i}{n} \log \frac{n_i}{n} \quad (2)$$

(3) The Shapelet with the largest information gain is selected as the best Shapelet.

#### 3.2 Dictionary Learning Algorithm

Dictionary learning is mainly used to reconstruct and reduce the dimensionality of the data in order to quickly find information about the data and classify the data [22]. The core idea of dictionary learning is to reconstruct an arbitrary signal by a linear combination of a specific set of dictionaries so that its coefficients are sparse.

A classical dictionary learning model representation is shown in equation (3):

$$\begin{aligned} \arg \min_{\alpha, Z} \frac{1}{2} \|T - \alpha Z\|_2^2 + \lambda \|\alpha\|_1 \\ s.t. \|Z_k\|_2^2 \leq c, k = 1, \dots, K \end{aligned} \quad (3)$$

where,  $T$  is the input dataset;  $Z = [Z_1, Z_2, \dots, Z_K] \in \mathbb{R}^{K \times m}$  is the dictionary to be learnt, where each element  $Z_K$  is an atom of the dictionary;  $\alpha \in \mathbb{R}^{n \times K}$  denotes the sparsity coefficients learnt from the data,  $\lambda$  is the weights for balancing the reconstruction error and sparsity; and  $c$  is the scaling factor to avoid dictionary distortions that produce a mundane solution.

Through the analysis of the Shapelet algorithm, it can be found that the processing of the time series is to carry out a linear combination of local shape features, which is similar to dictionary learning, and the combination of the two will greatly improve the efficiency and accuracy of classification. Therefore, the dictionary learning method is introduced into the Shapelet algorithm to improve the efficiency and accuracy of the Shapelet algorithm in processing the load curve.

#### 3.3 Construction of SDL-based industrial load recognition model

Through the analysis and optimization of industrial load, it can achieve the fine management of the industrial production process, improve the production efficiency and reduce the cost of energy consumption. The industrial load recognition based on Shapelet dictionary learning algorithm introduces the idea of dictionary to identify the temporal trajectory of the load curve. This approach improves both the computing speed and generalization ability of Shapelet as well as the efficiency and accuracy of load recognition. The model of SDL algorithm is shown in equation (4):

$$\begin{aligned} & \arg \min_{\alpha, q, d} \frac{1}{2} \sum_{i=1}^n \left\| T_i - \sum_{k=1}^K \alpha_{ik} Q(d_k, q_{ik}) \right\|_2^2 + \lambda \sum_{i=1}^n \|\alpha_i\|_1 \\ & s.t. \|d_k\|^2 \leq c, k=1, \dots, K \\ & 0 \leq q_{ik} \leq m-p, i=1, \dots, n; k=1, \dots, K \\ & \alpha_{ik} \geq 0, i=1, \dots, n; k=1, \dots, K \end{aligned} \quad (4)$$

where the input is a time series data set  $T$  and the outputs are: sparse coefficients  $\{\alpha_{ik}\}$ ; a Shapelet dictionary  $D = [d_1, d_2, \dots, d_k]$  consisting of  $K$  atoms, each atom  $d_k$  is a Shapelet of length  $p$ ;  $q_{ik}$  is a translation parameter. The hyperparameters to be tuned are: the number of Shapelets  $K$ , the length  $p$ , the weight  $\lambda$ , and the scaling factor  $c$ . For the sparse coefficients  $\alpha$  and the Shapelet dictionary  $D$  the initial values are given.

### Optimizing the Shapelet dictionary learning model

(1) Update the translation parameters

Fixing  $\alpha$  and  $d$  to find the optimal panning parameter  $q$ .

$$\begin{cases} \arg \min_{q_k} \frac{1}{2} \left\| T_i - \sum_{j \neq k} \alpha_j Q(d_j, q_j) - \alpha_k Q(d_k, q_k) \right\|_2^2 \\ s.t. \quad 0 \leq q_k \leq m-p, k=1, 2, \dots, K \end{cases} \quad (5)$$

Using the enumeration method, find the shortest Euclidean distance between the Shapelet and all subsequences to find the best  $q_k^*$ :

$$q_k^* = \arg \min_{q_k} \left\| d_k - \hat{t}^{[q_k+1, q_k+p]} \right\|^2 \quad (6)$$

where,  $\hat{t} = T_i - \sum_{j \neq k} \alpha_j Q(d_j, q_j)$  is the reconstructed residuals of all other Shapelets.

(2) Update the sparsity coefficients

Fix  $d$  and  $q$  to find the optimal sparsity coefficients  $\alpha$ .

$$\begin{cases} \arg \min_{\alpha_k} \frac{1}{2} \left\| T_i - \sum_{j \neq k} \alpha_j Q(d_j, q_j) - \alpha_k Q(d_k, q_k) \right\|_2^2 + \lambda |\alpha_k| \\ s.t. \quad \alpha_k \geq 0, k=1, 2, \dots, K \end{cases} \quad (7)$$

If there is no constraint, there exists an irreducible point of  $|\alpha_k|$  with respect to  $\alpha_k$ , whereas Eq. (7) has the property of being reducible in the domain of definition due to the non-

negative constraint  $\alpha_k$  of  $\alpha$ . The optimal  $\alpha_k$  can be found:

$$\alpha_k^* = \begin{cases} \frac{Q(d_k, q_k)^T \hat{t} - \lambda}{\|d_k\|^2} & / Q(d_k, q_k)^T \hat{t} > \lambda \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

where, a change in  $q_k$  changes the positivity and negativity of  $Q(d_k, q_k)^T \hat{t} - \lambda$ , which affects the derivability of Equation (8), so in each macrocycle, Steps 1 and 2 are iterated through a small cycle to reach convergence before being passed on to the next step.

(3) Update the Shapelet dictionary

Fixing  $\alpha$  and  $q$  to find the optimal Shapelet dictionary  $d$ .

$$d_k^* = \begin{cases} \sqrt{c} \hat{t} / \|\hat{t}\| & \text{if } \|\hat{t}\| / \sum_{i=1}^n \alpha_{ik}^2 > \sqrt{c} \\ \hat{t} / \sum_{i=1}^n \alpha_{ik}^2 & \text{otherwise} \end{cases} \quad (9)$$

in Eq.  $\tilde{t} = \sum_{i=1}^n \alpha_{ik} \hat{t}_i^{[q_{ik}+1, q_{ik}+p]}$

### Shapelet Transformation

Shapelet transformation operation is used to make the features available for classifier learning. The Shapelet transformation operation converts the time series into a feature space. After extracting all Shapelet dictionaries, the minimum distance criterion is used to determine the Euclidean distance between a Shapelet subsequence  $\bar{S}^k$  of length  $l$  and a subsequence  $X_i^l$  of the same length in the load curve.

### Construction of Random Forest Classification Model

Random Forest adopts the integration idea of Bagging algorithm and improves the generalization ability and stability of the overall model by taking putative regression sampling of the training dataset, then training an independent model for each subset, and finally combining the prediction results of these models to improve the generalization ability and stability of the overall model [23-24].

## 4. Analysis of examples

This paper presents a case study using the CEEMD-PCA-AE-Shapelet dictionary learning model to determine its effectiveness based on load data in an industrial park in Liaoning Province, China. All experiments were done on a server with Intel Xeon Gold 6248R CPU, NVIDIA Tesla V100 GPU, and Python 3.8, PyTorch 1.9.0, scikit-learn 1.0.2 software framework. The loads are typically 96-point industrial load data with 15-minute intervals, containing a total of 12,000 samples covering four types of loads: machinery, glass, steel, and cement. The training set is 80% data (9,600 entries), which is chronologically divided into the

first 70% for model training and the last 10% for validation set hyper-parameter tuning. The test set is 20% data (2,400 entries) covering different seasons and production cycles to ensure generalizability assessment. In practice, the raw load data first needs to be noise reduced to eliminate anomalous data. Finally, four electricity usage patterns such as machinery, glass, steel as well as cement are input into the Shapelet dictionary learning module as category labels and the best results are filtered by time and sample point comparisons. At the same time, in order to prove the effectiveness of the proposed method, the relevant comparisons are made in the examples. Table 1 shows the parameters of each module of the CEEMD-PCA-AE-Shapelet model.

Table 1. Parameters of the modules of the CEEMD-PCA-AE-Shapelet model

model	parameter name	parameter value
CEEMD-PCA-AE	CEEMD decomposition layer	8-layer IMF
	PCA retained variance ratio	95%
	AE network structure	encoder: 128→64 decoder: 64→128
	AE activation function	ReLU
Shapelet dictionary learning	Shapelet length range	5-30 samples
	number of dictionary atoms (K)	50
	sparsity weight ( $\lambda$ )	0.1
Random Forest classifier	number of decision trees	150
	maximum tree depth	20
	minimum number of leaf node samples	5

In outlier detection, recall and precision are more critical than accuracy, and the value of F1 not only includes recall, but also combines precision as a criterion for evaluation, and the higher the value, the more effective the detection method is. The AUC index is a key reference factor for evaluating the performance of the outlier detection model, and the value of AUC generally falls in the range of 0.5~1, and the value of AUC close to 1 means that the detection method is more effective. The AUC value generally falls within the range of 0.5~1, and the AUC value close to 1 means that the detection method is more effective.

Table 2. Comparison of anomaly detection results

data set	PCA	AE	CEEMD-PCA-AE	
machinery	AUC value	0.9367	0.9598	0.9813
	precision rate	0.91	0.93	0.89
	recall rate	0.93	0.95	0.9
	F1 value	0.92	0.94	0.9
glass	AUC value	0.9588	0.9682	0.9961
	precision rate	0.93	0.95	0.97
	recall rate	0.94	0.96	0.98

steel	F1 value	0.94	0.96	0.98
	AUC value	0.9505	0.9695	0.987
	precision rate	0.92	0.94	0.96
	recall rate	0.93	0.95	0.98
cement	F1 value	0.93	0.93	0.97
	AUC value	0.9428	0.9633	0.9813
	precision rate	0.81	0.83	0.89
	recall rate	0.83	0.85	0.9
	F1 value	0.83	0.85	0.9

Table 3. Multi-Algorithm Comprehensive Performance Comparison

method	accuracy (%)	F1 value	Training time (min)
PCA	82.3	0.81	5.2
AE	85.6	0.84	8.7
XGBoost	86.7	0.85	12.1
Wavelet-CNN	91.5	0.90	45.8
CEEMD-PCA-AE	95.8	0.94	18.7

As can be seen from Table 2, compared with using a single AE and PCA model, CEEMD-PCA-AE improves the AUC, precision, recall, and F1 value by about 3.13%, 5.11%, 4.46%, and 4.74%, respectively. As can be seen from Table 3, CEEMD-PCA-AE significantly outperforms AE, PCA, XGBoost [25], and Wavelet-CNN [26] models in terms of accuracy and F1 value, and the training time is better than the deep learning model, Wavelet-CNN model, and close to XGBoost. it can be seen that the method has a better anomaly detection capability.

Before the execution of the load recognition algorithm, the number of decision trees needs to be initialized. Figure 2 shows the variation of OBB estimation with respect to decision tree.

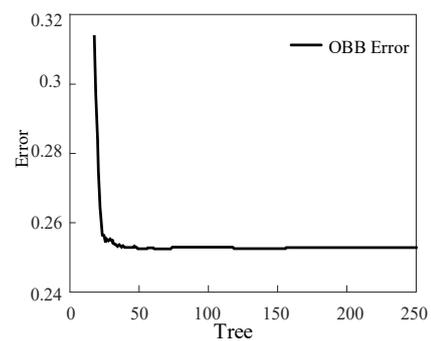


Figure 2. Plot of the change in the OBB estimate versus the decision tree

The generalization ability of the random forest model is closely related to the number of decision trees. Figure 3 demonstrates the trend of OBB error with the number of decision trees. The OBB error is the prediction error calculated by the random forest through the samples that are

not involved in the training of a single tree during the training process, which can effectively reflect the degree of model overfitting. When the number of trees in the forest is relatively small, the OBB error will have large fluctuations. With the gradual increase in the number of trees, the error of OBB begins to stabilize when the number of trees reaches 150, at which time the model error is reduced to a minimum. Therefore, when the model error is minimized, the corresponding values of trees are selected as the parameter inputs of the random forest. On this basis, the Gini coefficient formula and the solution method of optimal splitting value are used to determine the optimal splitting characteristics and the optimal splitting value of each splitting point until each tree grows completely. Based on the above SDL-RF to determine the load category, the results are shown below and Fig. 3 shows the load identification results.

Figure 3 illustrates the results of the SDL-RF model for identifying four categories of industrial loads. The curves for each category of loads are characterized as follows: category

1 shows periodic spikes; category 2 has a smooth baseline with intermittent high-power pulses. The blue curve in the figure shows the actual load data, and the red curve shows the model identification results. It can be seen that SDL-RF can accurately capture the timing characteristics of the load curves with a classification accuracy of 96.3% (see Table 1), which is significantly better than the traditional method. To further validate the effectiveness of the proposed method, accuracy metrics are compared with the results of Fully Convolutional Networks (FCN), Support Vector Machines (SVM), Random Forest (RF), and Shapelet Random Forest Classification Model (SRF). As shown in the classification result comparison in Figure 4, for the 50 data points represented by the blue dots, when the blue points are below the red line, the proposed method demonstrates higher accuracy. Compared to FCN, SVM, RF, and SRF, the SDL-RF method achieves higher accuracy for approximately 66%, 72%, 86%, and 58% of the total samples, respectively. The proposed SDL-RF algorithm exhibits superior classification accuracy.

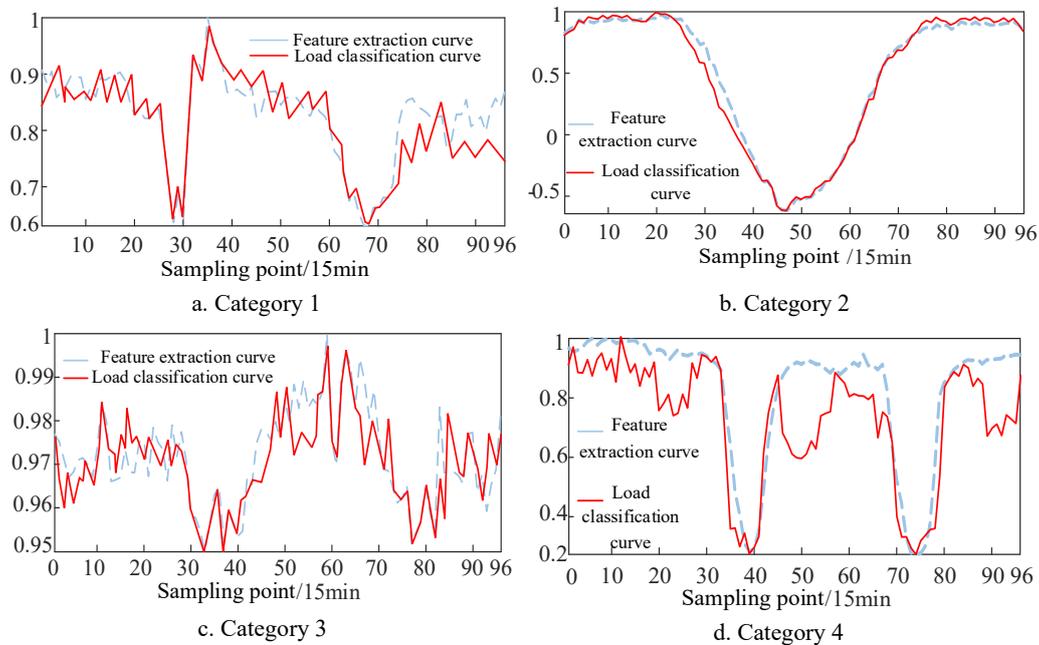
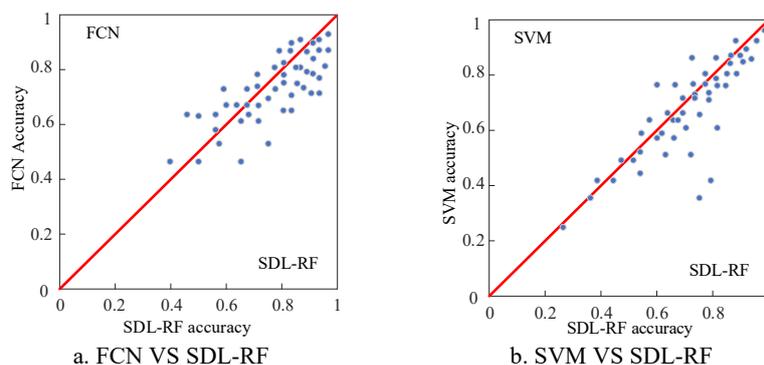


Figure 3. Load identification result chart



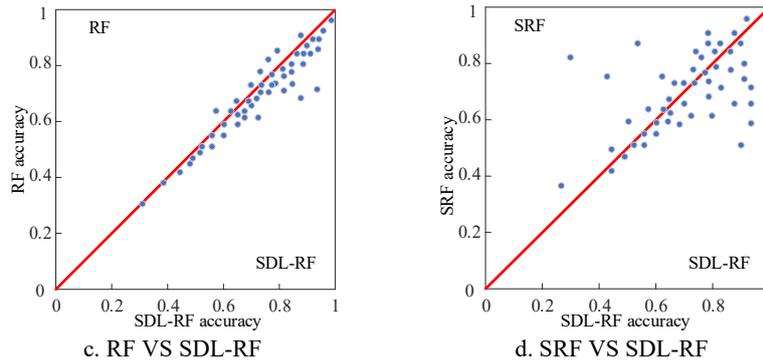


Figure 4. Comparison chart of classification results

To further validate the effectiveness of the algorithm, a comparison is made from the running time with RF and SRF. Figure 5 shows the running time comparison, from which it can be seen that when the data volume is more than 64MB, the running time of SRF and RF rises sharply, while that of SDL-R rises more gently. When the data volume is about 500MB, the SRF running time is about 3 times of SDL-RF, and the RF running time is about 2 times of DL-RF. Therefore, for massive load data, the SDL-RF method proposed in the paper has obvious superiority.

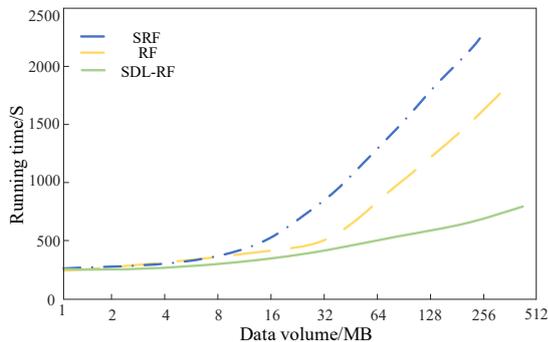


Figure 5. Runtime comparison chart

## 5. Conclusion

In this paper, a novel energy internet load identification method based on CEEMD-PCA-AE-Shapelet dictionary learning algorithm is proposed, which effectively solves the problems of sparsity, noise, and nonlinearity of industrial load data. Experimental results show that the model significantly improves the accuracy and efficiency of load identification and anomaly detection, and outperforms the traditional model in terms of accuracy, recall, and computational speed. The method has great potential to enhance energy management in industrial environments, contributing to more efficient and effective management of energy Internet loads. However, the method also has some limitations. First, although the model has been validated using data from Liaoning Province, its performance may vary when applied to different industrial

environments or geographic areas with different load patterns. In addition, the computational complexity of the model may hinder its real-time application, especially in systems with limited computational resources or small datasets. Future research could focus on improving the scalability of the model and reducing its computational cost to make it more suitable for real-time deployment in different environments. In addition, exploring the model's adaptability to different energy sectors could broaden its application scope and provide further insights into its global potential. Overall, the approach presented in this paper shows great potential for advancing energy management systems, but requires further optimization and testing.

## References

- [1] Manzoor A, Akram W, Judge M A, et al. Efficient economic energy scheduling in smart cities using distributed energy resources[J]. *Science and Technology for Energy Transition*, 2024, 79: 29.
- [2] Adeleye S A, Adebajji B, Awogbemi O. Renewable energy sources acceptability for decentralized energy system in Nigeria: Issues, challenges and prospects[J]. *Science and Technology for Energy Transition*, 2024, 79: 44.
- [3] Aziz M S, Ahmed S, Saleem U, et al. Wind-hybrid power generation systems using renewable energy sources-A review[J]. *International Journal of Renewable Energy Research*, 2017, 7(1): 111-127.
- [4] Wang S, Dang Q, Gao Z, et al. An innovative square root-untuned Kalman filtering strategy with full-parameter online identification for state of power evaluation of lithium-ion batteries[J]. *Journal of Energy Storage*, 2024, 104: 114555.
- [5] Wang S, Gao H, Takyi-Aninakwa P, et al. Improved multiple feature-electrochemical thermal coupling modeling of lithium-ion batteries at low-temperature with real-time coefficient correction[J]. *Protection and Control of Modern Power Systems*, 2024, 9(3): 157-173.
- [6] Williamson S S, Rimalapudi S C, Emadi A. Electrical modeling of renewable energy sources and energy storage devices[J]. *Journal of Power Electronics*, 2004, 4(2): 117-126.
- [7] Guo Z, Zhou K, Zhang C, et al. Residential electricity consumption behavior: Influencing factors, related theories and

- intervention strategies[J]. *Renewable and Sustainable Energy Reviews*, 2018, 81: 399-412.
- [8] Dong M, Sun M, Song D, et al. Real-time detection of wind power abnormal data based on semi-supervised learning Robust Random Cut Forest[J]. *Energy*, 2022, 257: 124761.
- [9] Wang X, Yao Z, Papaefthymiou M. A real-time electrical load forecasting and unsupervised anomaly detection framework[J]. *Applied Energy*, 2023, 330: 120279.
- [10] Chou J S, Telaga A S. Real-time detection of anomalous power consumption[J]. *Renewable and Sustainable Energy Reviews*, 2014, 33: 400-411.
- [11] Zhou K, Yang S. Understanding household energy consumption behavior: The contribution of energy big data analytics[J]. *Renewable and Sustainable Energy Reviews*, 2016, 56: 810-819.
- [12] Wang X, Ahn S H. Real-time prediction and anomaly detection of electrical load in a residential community[J]. *Applied Energy*, 2020, 259: 114145.
- [13] Cui M, Wang J, Yue M. Machine learning-based anomaly detection for load forecasting under cyberattacks[J]. *IEEE Transactions on Smart Grid*, 2019, 10(5): 5724-5734.
- [14] Rajabi A, Eskandari M, Ghadi M J, et al. A comparative study of clustering techniques for electrical load pattern segmentation[J]. *Renewable and Sustainable Energy Reviews*, 2020, 120: 109628.
- [15] Lei Y, Lin J, He Z, et al. A review on empirical mode decomposition in fault diagnosis of rotating machinery[J]. *Mechanical systems and signal processing*, 2013, 35(1-2): 108-126.
- [16] Wang S, Zhang S, Wen S, et al. An accurate state-of-charge estimation of lithium-ion batteries based on improved particle swarm optimization-adaptive square root cubature kalman filter[J]. *Journal of power sources*, 2024, 624: 235594.
- [17] Tanoni G, Principi E, Squartini S. Non-Intrusive Load Monitoring in industrial settings: A systematic review[J]. *Renewable and Sustainable Energy Reviews*, 2024, 202: 114703.
- [18] Ding Y, Chen Z, Zhang H, et al. A short-term wind power prediction model based on CEEMD and WOA-KELM[J]. *Renewable Energy*, 2022, 189: 188-198.
- [19] Yeh J R, Shieh J S, Huang N E. Complementary ensemble empirical mode decomposition: A novel noise enhanced data analysis method[J]. *Advances in adaptive data analysis*, 2010, 2(02): 135-156.
- [20] Hotelling H. Analysis of a complex of statistical variables into principal components[J]. *Journal of educational psychology*, 1933, 24(6): 417.
- [21] Rumelhart D E. Learning internal representations by error propagation[J]. *Parallel distributed processing: explorations in the microstructure of cognition*, 1986, 1: 319-362.
- [22] Rakotomamonjy A. Direct optimization of the dictionary learning problem[J]. *IEEE Transactions on Signal Processing*, 2013, 61(22): 5495-5506.
- [23] Agarwal S, Chowdary C R. A-Stacking and A-Bagging: Adaptive versions of ensemble learning algorithms for spoof fingerprint detection[J]. *Expert Systems with Applications*, 2020, 146: 113160.
- [24] Belgiu M, Drăguț L. Random forest in remote sensing: A review of applications and future directions[J]. *ISPRS journal of photogrammetry and remote sensing*, 2016, 114: 24-31.
- [25] Sagi O, Rokach L. Approximating XGBoost with an interpretable decision tree[J]. *Information sciences*, 2021, 572: 522-542.
- [26] Nguyen D C, Salamak M, Katunin A, et al. Vibration-based SHM of railway steel arch bridge with orbit-shaped image and wavelet-integrated CNN classification[J]. *Engineering Structures*, 2024, 315: 118431.