

A novel autoencoder approach with rescaling transformation distance metric for feature extraction to high-dimensional data

Jian Zheng¹, Hongchun Qu^{Corresp., 1, 2}, Zhaoni Li¹, Lin Li¹, Xiaoming Tang², Fei Guo²

¹ College of Computer Science and Technology, Chongqing University of Post and Telecommunications, Chongqing, China, China

² College of Automation,, Chongqing University of Posts and Telecommunications,, chongqing, China, China

Corresponding Author: Hongchun Qu
Email address: hcchyu@gmail.com

Feature extraction usually needs to rely on sufficient information of the input data, unfortunately, the distribution of the data upon a high-dimensional space is too sparse to provide sufficient information for feature extraction. Moreover, high dimensionality of the data also brings trouble for the searching of those features scattered in subspaces. As such, it is a tricky work for feature extraction from the data in a high-dimensional space. To address this issue, this paper proposes a novel autoencoder approach based on Mahalanobis distance metric of rescaling transformation. Through performing rescaling transformation on Mahalanobis distance metric, then the transformed Mahalanobis distance metric is introduced into the autoencoder, so as to improve the ability of feature extraction to the model. Results show that the proposed approach wins comparison methods in terms of both the accuracy of feature extraction and the linear separabilities of the extracted features. We indicate that distance metric-based methods are more suitable for extracting those features with linear separabilities from high-dimensional data than feature selection-based methods. In a high-dimensional space, evaluating feature similarity is relatively easier than evaluating feature importance, so that distance metric methods by evaluating feature similarity gain advantages over feature selection methods by assessing feature importance for feature extraction, while the latter is higher than the former in terms of computational efficiency.

A novel autoencoder approach with rescaling transformation distance metric for feature extraction to high-dimensional data

Jian Zheng ¹, Hongchun Qu^{1,2,*}, Zhaoni Li ¹, Lin Li¹, Xiaoming Tang ², Fei Guo²

¹ College of Computer Science and Technology, Chongqing University of Posts and Telecommunications, Chongqing, 400065, China

² College of Automation, Chongqing University of Posts and Telecommunications, Chongqing, 400065, China

Corresponding Author:

Hongchun Qu

Congwenlu Road, Chongqing, 400065, China

Email address: hcchu@gmail.com

Abstract

Feature extraction usually needs to rely on sufficient information of the input data, unfortunately, the distribution of the data upon a high-dimensional space is too sparse to provide sufficient information for feature extraction. Moreover, high dimensionality of the data also brings trouble for the searching of those features scattered in subspaces. As such, it is a tricky work for feature extraction from the data in a high-dimensional space. To address this issue, this paper proposes a novel autoencoder approach based on Mahalanobis distance metric of rescaling transformation. Through performing rescaling transformation on Mahalanobis distance metric, then the transformed Mahalanobis distance metric is introduced into the autoencoder, so as to improve the ability of feature extraction to the model. Results show that the proposed approach wins comparison methods in terms of both the accuracy of feature extraction and the linear separabilities of the extracted features. We indicate that distance metric-based methods are more suitable for extracting those features with linear separabilities from high-dimensional data than feature selection-based methods. In a high-dimensional space, evaluating feature similarity is relatively easier than evaluating feature importance, so that distance metric methods by evaluating feature similarity gain advantages over feature selection methods by assessing feature importance for feature extraction, while the latter is higher than the former in terms of computational efficiency.

Keyword: Autoencoder, Distance metric, Feature extraction

Introduction

High-dimensional data usually contains rich features, through extracting the important features, those irrelevant attributes in high-dimensional data can be filtered, thereby achieving data dimensionality reduction (Xue et al., 2015). Hence, feature extraction is considered to be one of the important methods for data dimension reduction (Bo et al., 2016)].

Feature extraction is a hot topic in recent years, aiming to gain the most valuable features from the input data (H.Tao et al., 2016; M. Luo et al., 2018). High dimensionality of data, the so-called the curse of dimensionality, brings negative effects for feature extraction (J.Gui et al., 2018; R. Chakraborty, N.R.Pal, 2015). Upon a low-dimensional space, those relations between the data are relatively compact but they may become sparse upon a high-dimensional space (Bing et al., 2021), e.g., the data space with more than 10 dimensionalities (A.-M., et al, 2011). Clearly, sparse relations between data are usually considered to be an unfavorable factor for feature extraction since feature extraction needs to rely on the relations between data (Bo et al., 2021). Beyond that, those latent features scattered in subspaces inside a high-dimensional space not only inspect the ability of methods to extract features (L. Wang et al., 2016), but also test their extraction efficiency. Hence, it is a challenge for feature extraction from high-dimensional data.

Recently, some opinions have been proposed for feature extraction, for instance, distance metric-based methods, where, the typical representative is the well-known Mahalanobis distance-based methods, which evaluates the similarity between samples using the covariance matrix of data (R. De et al, 2000). Furthermore, S.-Y (2018) et al proposed the intrinsic semi-supervised metric learning (ISSML) based on a distance metric for feature extraction. Similar, the methods implemented in (P. Zadeh et al, 2016) and (H. Luo, 2017) also applied distance metrics. Certainly, also including, the information-theoretic metric learning is (ITML) (J. Mei et al, 2014) employed a distance metric to obtain features. These methods in (S. Ying et al, 2018; P. Zadeh et al, 2016; H. Luo, 2017; J. Mei et al, 2014) address the issues of symmetric positive-definite matrix minimization during feature extraction, but there are several problems in them, 1) since most of them use iterative calculation while performing feature selection, optimization issues have to be addressed iteratively. 2) Most of them need to rely on parameter selection to obtain those desired features. Usually, feature selection-based methods are also considered to be used for feature extraction. Such methods achieve feature extraction through analyzing the information of feature subsets, for example, the cheap feature selection method based on k -means algorithm (Marco et al, 2021) selects the m features with the highest relevance measure through obtaining a clustering for each subset of features. Although the method (Marco et al, 2021) is a novel measurement for feature relevance, which is beneficial for feature selection, however, calculating per subset of features needs to spend a lot time cost. In order to reduce the correlation between features, some measurements for quickly assessing features are proposed, e.g., the information entropy metric (T. X. et al, 2019), whereas the method (T. X. et al, 2019) has a bias toward features, which may result in appearing selecting deviation during feature extraction. Another kind of feature selection method depends on eigen decomposition, such as, locally linear embedding (LLE) (R. Hettiarachchi, J. F. Peters, 2015; Ugochukwu Ejike Akpudo,

Jang-Wook Hur, 2020), multi-manifold discriminant isometric feature mapping (MMD-ISOMAP) (*Bo.Y et al, 2016*), ISOMAP-KL (*Alaor Cervati Neto, Alexandre L. M. Levada, 2020*), however, they cannot assess the importance of the features in the background space explicitly.

Neural network-based methods are favored because of excellent feature capture ability (*Hong.C et al, 2022*), e.g., Multilayer Perceptron Neural Network (*K. Sun et al., 2017*). For dimension reduction, feature extraction and data compression, autoencoder-based networks provide an interpretable approach for the unknown meaningful insights (*Ang et al., 2017*) by learning non-identity mapping functions (*Jian et al., 2022*), for instance, *Rami et al (2022)* developed interpretable data representation for data dimensionality reduction using Logic-Oriented and Granular Logic Autoencoders, and such as, Autoencoder (*Angshul Majumdar, 2019*) for image compression, and Blind Denoising autoencoder (*Fei et al, 2020*) for denoising. In addition, sparse autoencoders are used as an unsupervised feature extractor to serve data dimensionality reduction, feature extraction and data mining (*Z.Qiang et al., 2018*), e.g., *K.-J et al (2018)* proposed Sparse Autoencoder (SAE) for feature extraction of ferroresonance overvoltage waveforms in power distribution systems.

In this study, our motivation is to extract the features with linear separabilities from the data in a high-dimensional space. Thus, we proposed a novel autoencoder method based on Mahalanobis distance metric of rescaling transformation. The proposed method does not have to address any optimization issue, and also it can focus on the whole data distribution.

We summarize the main contributions of this work as follows:

(i) Distance metric-based methods are more suitable for extracting those features with linear separabilities from high-dimensional data than feature selection-based methods.

(ii) Assessing feature similarity in a high-dimensional space is relatively easier than evaluating feature importance, therefore, distance metric approaches by evaluating feature similarity have more advantages than feature selection approaches by evaluating feature importance in terms of feature extraction.

(iii) The computational time of distance metric-based algorithms is higher than that of feature selection-based algorithms upon a high-dimensional space.

This paper is organized as follows. Section 2 describes the proposed method and implements the proposed model, including training for the model and parameter configuration. Experiment datasets, competing methods, and experiment description are given in Section 3. Section 4 presents experiment results. Section 5 draws conclusions.

Methods

Theory

Given a sample $X = \{x_i | 1 \leq i \leq N\}$, and $X \subseteq \mathbb{R}^d$. \mathbb{R}^d is the d -dimensional Euclidean space. P is the probability distribution of X , denoted as original probability distribution. $\mu(X)$ and Γ_x are the mean vector and the covariance matrix of X , respectively. Let us assume that $Z = \{z_j | 1 \leq j \leq N\}$ is the reconstructed X , and $Z \subseteq \mathbb{R}^d$. Q is the probability distribution of Z , denoted as approximate

probability distribution. Similar, $u(Z)$ and Γ_z are the mean vector and the covariance matrix of Z , respectively. The K-L divergence (D. Tao et al, 2009) between the two distributions P and Q is given in Eq. (1).

$$K(P \parallel Q) = \frac{1}{2} [\log |\Gamma_z| - \log |\Gamma_x| + \text{tr}(\Gamma_z^{-1} \Gamma_x) + \text{tr}(\Gamma_z^{-1} D_{xz})] \quad (1)$$

Where $|\Gamma| = \det(\Gamma)$. The $\text{tr}(\bullet)$ is the trace of a matrix. $D_{xz} = (u(X) - u(Z))(u(X) - u(Z))^T$ is a symmetrical matrix. Training a distance metric is equivalent to finding a rescaling of a sample which replaces each x_i with $M^T x_i$ (Lin et al, 2019), so the K-L divergence in Eq. (1) can be converted into Eq. (2), having

$$K_L^*(P \parallel Q) = \frac{1}{2} [\log |M^T \Gamma_z M| - \log |M^T \Gamma_x M| + \text{tr}((M^T \Gamma_z M)^{-1} (M^T (\Gamma_x + D_{xz}) M))] \quad (2)$$

Where M is a metric matrix and satisfies $A^* = MM^T$, and $M \in \mathbb{R}^{d \times d_0}$, $d_0 \leq d$. The K-L divergence in Eq. (2) is the rescaling transformation for the K-L divergence in Eq. (1) using the distance metric matrix A^* . To reduce the difference between the approximate distribution Q and the original distribution P , we consider Mahalanobis distance metric for K-L divergence in Eq. (2), having

$$K-L(\mathcal{d}_{A^*}) = K_L^*(P \parallel Q) + \sum_{1 \leq i, j} d_{A^*}(x_i, z_j) \quad (3)$$

$d_{A^*}(x_i, z_j)$ is Mahalanobis distance between x_i and z_j using A^* . The advantage of doing this is that the Mahalanobis distance using A^* can appropriately measure similarities between the input sample and the reconstructed input sample because of non-negativity (i.e., $d_{A^*}(x_i, z_j) \geq 0$), distinguishability (i.e., $d_{A^*}(x_i, z_j) = 0 \Leftrightarrow x_i = z_j$) and symmetry (i.e., $d_{A^*}(x_i, z_j) = d_{A^*}(z_j, x_i)$) (Lin et al, 2019). Eq. (4) gives the calculation of $d_{A^*}(x_i, z_j)$, where A^* can be decomposed as $A^* = MM^T$.

$$d_{A^*}(x_i, z_j) = \sqrt{(x_i - z_j)^T A^* (x_i - z_j)} \quad (4)$$

Model implementation

A classic auto encoder (AE) consists of an input layer, a hidden-layer and an output layer. For AE, the loss error is often measured by using the distance between the original input instance, the predicted instances, and the reconstructed instance (L. Theis et al, 2017). Usually, using divergence metrics or expanding autoencoder structures (e.g., enlarging the number of hidden layers) is more helpful for autoencoders to characterize the data distribution and to learn the desired representations (Weining et al, 2017). As such, we designed an autoencoder with multiple-hidden layers, namely m-AE, and $m \geq 1$, as shown in Fig.1. In addition, the K-L divergence in Eq. (3) was used to increase the ability of m-AE to capture low-dimensional feature representations. The loss error $\nabla_{KL}(\mathbf{w}, \mathbf{b})$ in m-AE is given as follows:

$$\nabla_{KL}(\mathbf{w}, \mathbf{b}) = \sum \|e_X - e_Z\|^2 + K-L(\mathcal{d}_{A^*}) \quad (5)$$

Where e_X , e_Z are the inputting and the reconstructed inputting, respectively. $\nabla_{KL}(\mathbf{w}, \mathbf{b})$ is updated through using the backpropagation manner.

To better train the proposed model, we carefully studied part hyper parameters in the model. For the rest of hyper parameters, their default values were used.

(i) Optimizer. Common optimizers are Adam, RMSprop, SGD, Momentum, Nesterov, etc. However, we selected Adam as the optimizer of m-AE, since Adam has the ability to handle sparse gradients (Diederik P. Kingma, Jimmy Lei Ba, 2015). Compared with other optimizers, Adam is more suitable for high-dimensional data. Moreover, Adam can provide different adaptive learning rates for different hyper parameters.

(ii) Activation function. Gradient vanishing is easily to be induced during passing gradients backwards for neural networks, in this case, the probability of gradient vanishing caused by activation function Sigmoid is relatively high. Similar to Sigmoid, activation function tanh also suffers from this problem. While for activation function ReLu, the phenomenon of gradient vanishing is partially alleviated, meaning that gradient vanishing does not appear in the positive interval of ReLu. Furthermore, ReLu converges much faster than Sigmoid and Tanh. Therefore, we chose ReLu as the activation function of m-AE.

(iii) Iteration epoch. We dynamically adjust the iteration epoch according to training accuracy. For instance, when training accuracy starts to change from large to small, we reduce iteration epoch in order to prevent over-fitting. When the difference in accuracy between training and testing is minimal, the current iteration epoch can be accepted and training procedure is stopped.

We give the training algorithm for m-AE in Algorithm 1. In the algorithm, the training set $Train_set$ is divided into two datasets $TCro_train$, $TCro_val$ in step 1. Since m-AE has multiple hidden layers, we set m in the range of O_m , in order to determine the m , the dataset $TCro_train$ is used to train m-AE. The data set $TCro_val$ is used for the validation of the network structure of m-AE. To get the optimal m , denoted as m_{opt} , the cross-validation is implemented in step 2 to step 18, where the procedure of step 6 and step 10 describes the calculation process of loss error $\nabla_{KL}(\mathbf{w}, \mathbf{b})$. After gaining the optimal m , m-AE is trained using the training set $Train_set$. Using backpropagation manner updates network parameters until m-AE can converge, as shown in step 18 to step 28. The procedure shown in step 29 to step 33 indicates that the maximum training accuracy $Train_acc$ are outputted and the well trained m-AE is saved.

Algorithm 1. Training for m-AE.

Input: Training set $Train_set$, $A^* = I \in \mathbb{R}^{d \times d}$ is an identity matrix, iteration epoch T , L , parameter O_m .

Output: Training accuracy $Train_acc$.

Begin:

1. $Train_set$ is divided into $TCro_train$, $TCro_val$;

2. **for** $t=1$ to T **do**:

3. **foreach** m **in** O_m :

4. Decompose A^* as satisfying $A^* = MM^T$ using eigen decomposition.

5. Calculate loss error $\nabla_{KL}(\mathbf{w}, \mathbf{b})$ using Eq. (5) and the procedure is summarized as following:

6. The procedure:

7. Calculate $K_L^*(P \parallel Q)$ using Eq. (2).

8. Calculate $d_{A^*}(x_i, z_j)$ using Eq. (4).

9. Take $K_L^*(P \parallel Q)$ and $d_{A^*}(x_i, z_j)$ into Eq. (3) to calculate $K-L(\mathcal{D}_{A^*})$.

10. For any x_i, x_j , calculate $\nabla_{KL}(\mathbf{w}, \mathbf{b})$ using Eq. (5).

11. Calculate training accuracy $T_acc = \text{m-AE}(T^{Cro_train}, m; t)$;

12. Validate m-AE using data set T^{Cro_val} ;

13. Calculate validation accuracy $V_acc = \text{m-AE}(T^{Cro_val}, m; t)$

14. Update weight $\mathbf{w} \leftarrow \mathbf{w} + \nabla \mathbf{w}$.

15. Update A^* as MM^T .

16. Until A^* and hyper parameters converge.

17. **end foreach**

18. **end for**

19. Get the optimal value of m , i.e., $m_{opt} = \arg \max(V_acc)$;

20. **for** $l=1$ **to** L **do**:

21. Decompose A^* as satisfying $A^* = MM^T$ using eigen decomposition.

22. Train m-AE using training set $Train_set$ and m_{opt} ;

23. Update network parameters using the optimizer Adam;

24. Calculate loss error $\nabla_{KL}(\mathbf{w}, \mathbf{b})$ using Eq. (5);

25. Calculate training accuracy $Training_acc(l) = \text{m-AE}(Train_set, m_{opt})$;

26. Update A^* as MM^T ;

27. Using backpropagation manner updates network parameters;

28. **end for**

29. Select the l so that $l_{max} = \arg \max(Training_acc(l))$;

30. Get the maximum training accuracy $Train_acc$ in the l_{max} -th iteration;

31. $Train_acc = \text{m-AE}(Train_set, m_{opt}, l_{max})$;

32. Output $Train_acc$

33. Save the well trained m-AE($Train_set, m_{opt}, l_{max}$);

End

EXPERMENTS

Datasets and assessment metrics

To verify the performance of the proposed m-AE, we selected 4 benchmark datasets with different data dimensions from the UCI machine learning repository (C. L. Blake, C. J. Merz, 1998). The attributes of the 4 benchmark datasets are summarized in Table 1.

Receiver operating characteristic curve (ROC) and corresponding area under curve (AUC) are usually used to assess the precision of machine learning methods. Therefore, AUC is taken as the assessment metric of method precision.

Competing and benchmark methods

Since m-AE applies the distance metric of rescaling transformation, the methods based on a distance metric were used for comparisons, including ISSML (*S.-Y et al, 2018*) and ITML (*J. Mei et al, 2014*). Certainly, the method based on feature selection was also considered, i.e., MMD-ISOMAP (*Bo.Y et al, 2016*). In addition, autoencoder-based approaches were used as a comparison, e.g., the SAE (*K.-J et al, 2018*). Furthermore, to further examine the effects of the distance metric of rescaling transformation on the performance of m-AE, a benchmark model was developed with m-AE as a reference. The developed benchmark model used the same structure and parameter configuration of m-AE without using the distance metric of rescaling transformation, namely AE-BK.

We implemented the corresponding algorithms of the six models using Python on Tensorflow framework. While for those parameters of competing methods, we adopted those values observed in the corresponding literature. Certainly, unless otherwise stated, the five corresponding algorithms all run on the same GPU and apply the same experimental configuration settings.

Experiment description

Experiments were conducted on the four benchmark datasets in order to validate the ability of these six models to extract features and their efficiency.

Experiment I. To test the robustness of m-AE. The proposed m-AE has multiple hidden layers, since the number of hidden layers (i.e., the m) significantly affects the precision of feature extraction, the m needs to be firstly verified, i.e., robustness testing of the model, let m set in the range of $\{1, 2, 3, 4, 5, 7, 10, 15, 20\}$.

Experiment II. To test the ability of feature extraction for the six models. The six models were run on the four benchmark datasets, and then the testing results were analyzed.

Experiment III. To compare the efficiency of our method with competing methods. These methods were performed on four benchmark datasets and observed their running time.

Ablation experiments. To verify that using the distance metric of rescaling transformation can be beneficial for extracting linearly separable features, the ablation experiments were also designed.

In addition, to eliminate randomness during the experiment and present an objective result, we used cross-validation to verify the six models. We randomly selected two datasets from the four benchmark datasets as the training set to train the six models. Once the six models were well trained, they were tested on the four benchmark datasets, respectively. The process was repeated five times, independently, then we took the average of five testing results as a measurement.

RESULTS

Experiments on robustness

Results in Fig.2 and Fig.3 show that the performance of the proposed m-AE and the benchmark model AE-BK improves along with increasing of m , and then the performance remains stable when m reaches a certain size, i.e., $m=3$. This means that m-AE and AE-BK are not sensitive to large m on the four benchmark datasets, i.e., their network structures are robust within a reasonable range. Therefore, let m be equal to 3 in subsequent experiments.

Comparisons of accuracy extraction

Results in Table 2 show that the proposed m-AE wins the four competing models and the benchmark model in the accuracy of feature extraction on all considered instances. For competitors, ISSML, ITML and SAE outperform MMD-ISOMAP in most benchmark datasets for the extracted accuracy.

Comparisons of linear separability.

The results of ablation experiments in Fig.4 show that compared with the models without using distance metrics, e.g., AE-BK, SAE, the models using distance metrics (including m-AE, ISSML, ITML) perform much better on most datasets in the extracted accuracy of the features with linear separabilities. Similar, the models using distance metrics also win the model using feature selection, as shown in Fig.5. To observe the linear separabilities of the extracted features from the four benchmark datasets, we projected these extracted features onto 2-dimensional space, and then visualized them. Fig.6 displays the results of visualized distribution on the four benchmark datasets by the six models. The visualized results show that it is optimal for the separation distance between different types of features extracted by m-AE, meaning that compared with competing and benchmark models, m-AE is a winner in terms of the linear separabilities of the extracted features. Together, these results imply that distance metric-based methods have advantages over feature selection-based methods in terms of extracting the features with linear separabilities.

Running time

Fig.7 displays the running time of methods. Obviously, the advantage of m-AE in running time is not as significant as that in both the extracted accuracy and the linear separabilities of the extracted features. MMD-ISOMAP spends less in running time on most benchmark datasets than distance metric-based methods, meaning that the execution efficiency of feature selection-based methods is higher than that of distance metric-based methods when running upon a high-dimensional space. Distance metric-based methods take a lot of time to calculate the distance between each point pair upon a high-dimensional space, so as to increase the running time.

Discussion

Compared with the competitors, the proposed m-AE has outstanding advantage in term of both

the accuracy of feature extraction and the linear separabilities of the extracted features on high-dimensional data. We interpret it as following. On one hand, Mahalanobis distance in Eq. (3) can appropriately measure similarities between the input sample and the reconstructed input sample, so as to minimize the loss error of m-AE in Eq. (5). As such, m-AE gains the desired accuracy of feature extraction. On the other hand, we performed a rescaling on K-L divergence metric in Eq. (2) by using A^* in Eq. (4), which effectively allows the extracted features to present linear separabilities, because the rescaling can maximized the classification distance between the extracted different types of features. Hence, the features extracted by m-AE present linear separabilities than competitors. Overall, m-AE outperforms the competitors in extracted accuracy and the linear separabilities of the extracted features.

In a high-dimensional space, distance metric-based methods easily evaluates the feature similarity by calculating the distance between the data, however, feature selection-based methods relatively difficulty assess the feature importance. Therefore, distance metric-based methods, e.g., ISSML (*S.-Y et al, 2018*) and ITML (*J. Mei et al, 2014*), are more suitable for extracting those low-dimensional features with the linear separability from high-dimensional data than feature selection-based methods. However, the computational time of feature selection-based methods, e.g., MMD-ISOMAP (*Bo.Y et al, 2016*), is lower than that of distance metric-based methods in a high-dimensional space, since distance metric-based methods spend too much in calculating the distance between each point pair.

Although autoencoders have excellent feature capture capabilities, they may perform poorly in extracting linearly separable features, e.g., SAE (*K.-J et al, 2018*). Whereas, this deficiency of autoencoders can be remedied by introducing a distance metric. Certainly, there are many methods of distance metrics, e.g., Wasserstein distance metric (*Na et al, 2019; Jian et al, 2022*), Bhattacharyya distance metric (*Mariucci E, Reiß M, 2017*).

Conclusions

This paper proposed a novel autoencoder method based on Mahalanobis distance metric of rescaling transformation to extract low-dimensional features from the data in the high-dimensional space. Results on show that compared with competing methods, the proposed method is a winner in both the accuracy of feature extraction and the linear separabilities of the extracted features. We find that the linear separabilities of those features obtained by the distance metric-based methods are better than that of obtained by the feature selection-based methods. Compared with feature selection-based methods, distance metric-based methods have more advantages for low-dimensional feature extraction, however, feature selection-based methods are higher than that of distance metric-based methods in terms of computational efficiency.

Funding

The funding received from the National Natural Science Foundation of China under grant #61871061

Data availability

The datasets in this work can be found at <http://archive.ics.uci.edu/ml/>

Competing interests

The authors declare no conflict of interest.

References

- Alaor Cervati Neto, Alexandre L. M. Levada. 2020. ISOMAP-KL: a parametric approach for unsupervised metric learning. *2020 33rd SIBGRAPI Conference on Graphics, Patterns and Images (SIBGRAPI)*. 1-1.
- A.-M. Zhou, K. D. Kumar, Z.-G. Hou. 2011. Finite-time attitude tracking control for spacecraft using terminal sliding mode and Chebyshev neural network. *IEEE Trans. Syst., Man, Cybern. B, Cybern.* **41**(4):950-963.
- Ang Jun Chin, Andri Mirzal, Habibollah Haron. 2016. Supervised, Unsupervised and Semi-supervised Feature Selection: A Review on Gene Selection. *IEEE Transactions on Computational Biology and Bioinformatics* **13**(5):971-989
- Angshul Majumdar. 2019. Blind Denoising Autoencoder. *IEEE Transactions on Neural Networks and Learning Systems* **30**(1):312-317.
- Bing Tu, Xiaolong Liao, Chengle Zhou. 2021. Feature Extraction Using Multitask Superpixel Auxiliary Learning for Hyperspectral Classification. *IEEE Transactions on Instrumentation and Measurement* **70**:1-16.
- Bo Peng, Shuting Wan, Ying Bi. 2021. Automatic Feature Extraction and Construction Using Genetic Programming for Rotating Machinery Fault Diagnosis. *IEEE Transactions on Cybernetics* **51**(10):4909-4923
- Bo Tang, Steven Kay, Haibo He. 2016. Toward Optimal Feature Selection in Naive Bayes for Text Categorization. *IEEE Transactions on Knowledge and Data Engineering* **28**(9):2508-2521.
- Bo. Yang, M. Xiang, Y. Zhang. 2016. Multi-manifold discriminant isomap for visualization and classification. *Pattern Recognition* **55**:215-230.
- C. L. Blake and C. J. Merz. 1998. UCI Repository of Machine Learning Databases, Department of Information and Computer Science. [Online]. Available: <http://www.ics.uci.edu/ml/learn/MLRepository.html>
- Diederik P. Kingma, Jimmy Lei Ba. 2015. Adam: A method for stochastic optimization. arXiv preprint. arXiv:1412.6980v8.
- D. Tao, X. Li, X. Wu, S. J. Maybank. 2009. Geometric mean for subspace selection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **31**(2):260-274.
- Fei Yang, Luis Herranz, Joost van de Weijer. 2020. Variable Rate Deep Image Compression With Modulated Autoencoder. *IEEE Signal Processing Letters* **27**:331-335.

- 389 **H. Luo. 2017.** Hyperspectral image classification using metric learning in one-dimensional
390 embedding framework. *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.* **10(5)**:1987-
391 2001.
- 392 **Hongchun Qu, Jian Zheng, Xiaoming Tang. 2022.** Effects of loss function and data sparsity
393 on smooth manifold extraction with deep model. *Expert Systems With Applications* **198**:1-
394 10.
- 395 **H. Tao, C. Hou, F. Nie. 2016.** Effective Discriminative Feature Selection With Nontrivial
396 Solution. *IEEE Transactions Neural Network Learning System* **27(4)**:796-808.
- 397 **J. Gui, Z. Sun, S. Ji. 2018.** Feature Selection Based on Structured Sparsity: A Comprehensive
398 Study. *IEEE Transactions Neural Network Learning System* **28 (7)**:1490-1507.
- 399 **Jian Zheng, Hongchun Qu, Zhaoni Li, Lin Li, Xiaoming Tang. 2022.** An irrelevant attributes
400 resistant approach to anomaly detection in high-dimensional space using a deep hyper
401 sphere structure. *Applied Soft Computing* **116** :1-20.
- 402 **J. Mei, M. Liu, H. R. Karimi. 2014.** Logdet divergence-based metric learning with triplet
403 constraints and its applications. *IEEE Transaction Image Process.* **23(11)**:4920-4931.
- 404 **K. Sun, S. H. Huang, D. S. Wong. 2017.** Design and Application of a Variable Selection
405 Method for Multilayer Perceptron Neural Network With LASSO. *IEEE Transactions*
406 *Neural Networks Learning System* **28(6)**:1386-1396.
- 407 **Kunjin Chen, Jun Hu, Jinliang He. 2018.** A Framework for Automatically Extracting
408 Overvoltage Features Based on Sparse Autoencoder. *IEEE Transactions on Smart Grid*
409 **9(2)**:594-604.
- 410 **Lin Feng, Huibing Wang, Bo Jin. 2019.** Learning a Distance Metric by Balancing KL-
411 Divergence for Imbalanced Datasets. *IEEE Transactions on Systems, Man, and*
412 *Cybernetics: Systems* **49(12)**:2384-2395.
- 413 **L. Theis, W. Shi, A. Cunningham. 2017.** Lossy image compression with compressive
414 autoencoders. *In Proc. Int. Conf. Learn. Representations*.1-19.
- 415 **L. Wang, Y. Wang, Q. Chang. 2016.** Feature selection methods for big data bioinformatics: A
416 survey from the search perspective. *Methods* **111**:21-31.
- 417 **Marco Capó, Aritz Pérez, Jose A. Lozano. 2021.** A Cheap Feature Selection Approach for the
418 K-Means Algorithm. *IEEE Transactions on Neural Networks and Learning Systems* **32(5)**:
419 2195-2208
- 420 **Mariucci E, Reiß M. 2017.** Wasserstein and total variation distance between marginals of Levy
421 processes. *ArXiv*:1710.02715.
- 422 **M. Luo, F. Nie, X. Chang. 2018.** Adaptive Unsupervised Feature Selection With Structure
423 Regularization. *IEEE Transactions Neural Network. Learning System* **29(4)**:944-956.
- 424 **Na Lei, Kehua Su, Li Cui. 2019.** A geometric view of optimal transportation and generative
425 model. *Computer Aided Geometric Design* **68**:1-28.
- 426 **P. Zadeh, R. Hosseini, S. Sra. 2016.** Geometric mean metric learning. *In Proc. Int. Conf. Mach.*
427 *Learning*. 2464- 2471.

- 428 **R. Hettiarachchi, J. F. Peters. 2015.** Multi-manifold lle learning in pattern recognition. *Pattern*
429 *Recognition* **48(9)**:2947-2960.
- 430 **Rami Al-Hmouz, Witold Pedrycz, Abdullah Balamash. 2022.** Logic-Oriented Autoencoders
431 and Granular Logic Autoencoders: Developing Interpretable Data Representation. *IEEE*
432 *Transactions on Fuzzy Systems* **30(3)**:869-877.
- 433 **R. Chakraborty, N. R. Pal. 2015.** Feature Selection Using a Neural Framework With
434 Controlled Redundancy. *IEEE Transactions Neural Network Learning System* **26(1)**:35-50.
- 435 **R. De Maesschalck, D. Jouan-Rimbaud, D. L. Massart. 2000.** The Mahalanobis distance,
436 Chemometrics Intell. Lab. System **50(1)**:1-18.
- 437 **S. Ying, Z. Wen, J. Shi. 2018.** Manifold preserving: An intrinsic approach for semisupervised
438 distance metric learning. *IEEE Transaction. Neural Networks Learning System* **29(7)**:
439 2731-2742.
- 440 **T. X. Pham, P. Siarry, H. Oulhadj. 2019.** A multi-objective optimization approach for brain
441 MRI segmentation using fuzzy entropy clustering and region-based active contour methods.
442 *Magn. Reson. Image* **61**:41-65.
- 443 **Ugochukwu Ejike Akpudo,Jang-Wook Hur. 2020.** Intelligent Solenoid Pump Fault Detection
444 based on MFCC Features, LLE and SVM. *2020 International Conference on Artificial*
445 *Intelligence in Information and Communication (ICAIIIC)*. 1-3.
- 446 **Weining Lu, Yu Cheng, Cao Xiao. 2017.** Unsupervised Sequential Outlier Detection with Deep
447 Architectures. *IEEE Transactions on Image Processing* **26(9)**:4321-4330.
- 448 **Xue, B., Zhang, M., Browne. 2015.** A Survey on Evolutionary Computation Approaches to
449 Feature Selection. *IEEE Transactions on Evolutionary Computation* **20(4)**:606-626.
- 450 **Zhiqiang Wan, Haibo He, Bo Tang. 2018.** A Generative Model for Sparse Hyperparameter
451 Determination. *IEEE Transactions on Big Data* **4(1)**:2-10.

Figure 1

The structure of the proposed m-AE

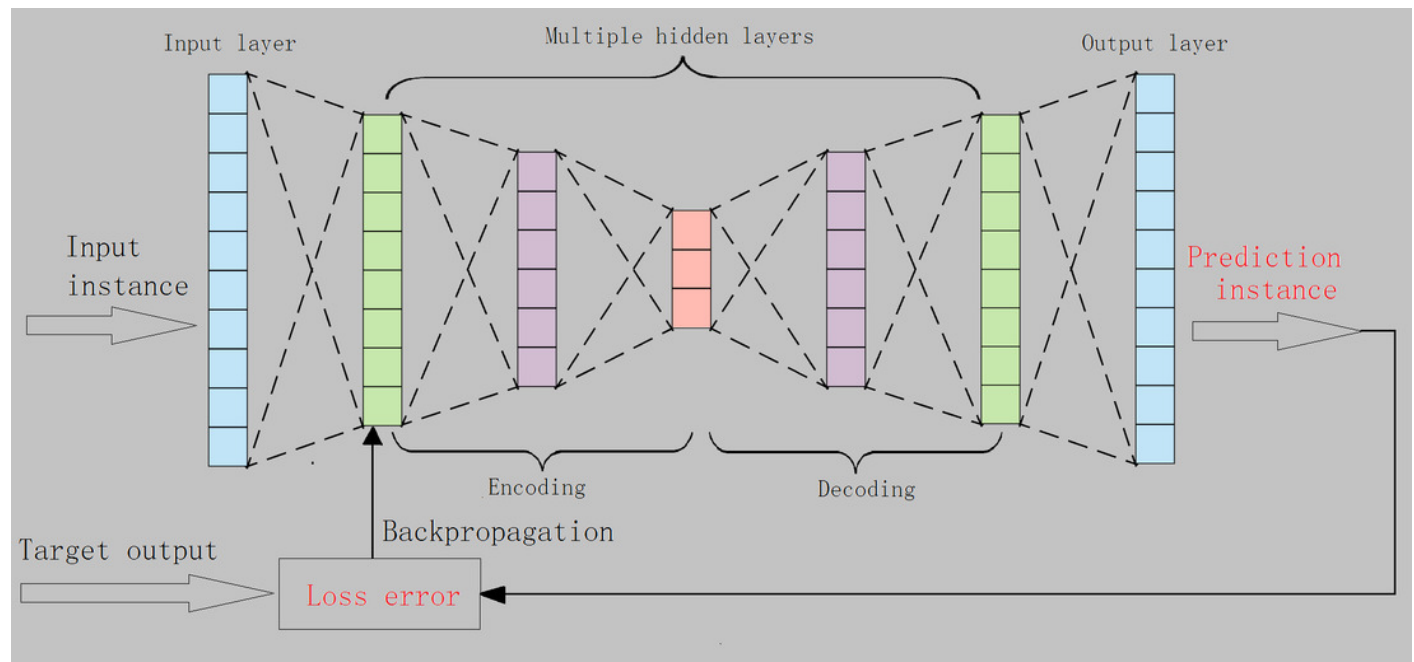


Figure 2

Validation of robustness for the proposed m-AE

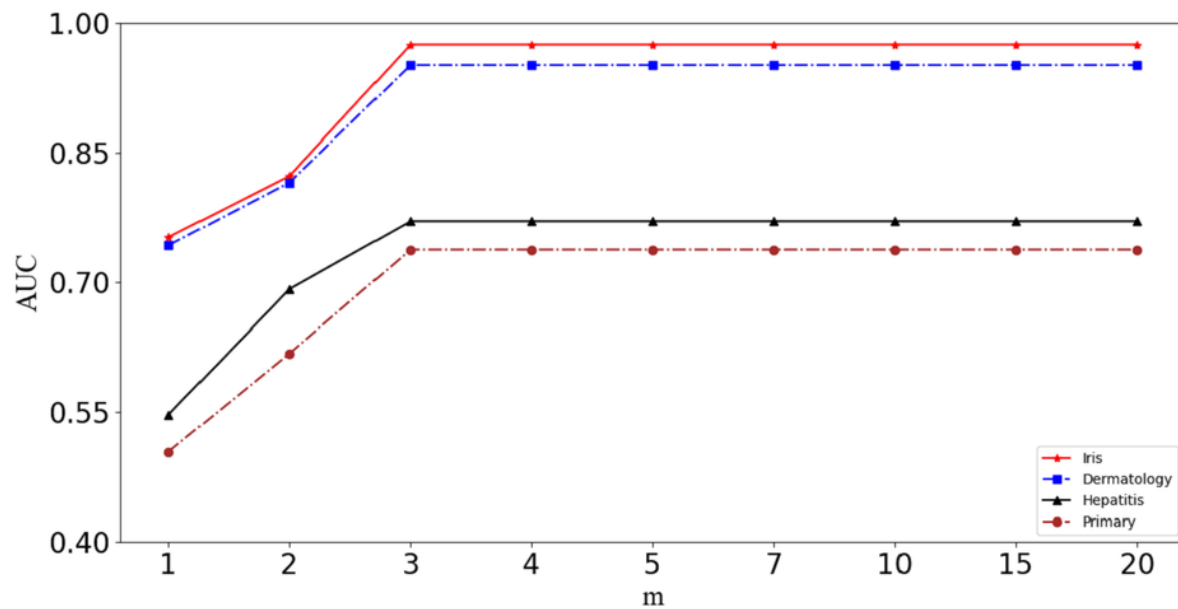


Figure 3

Validation of robustness for benchmark model AE-BK

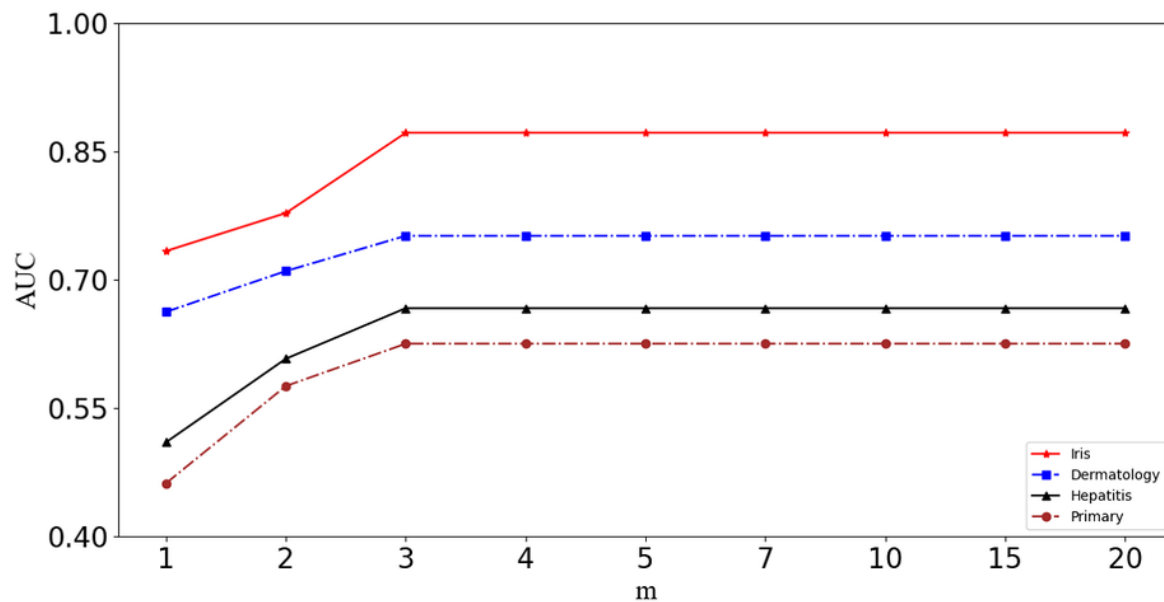


Figure 4

Results of ablation experiments

Comparisons between with distance metrics and without distance metrics. The models with distance metrics are marked as the symbol \checkmark . The models without both distance metrics and feature selection are marked as the symbol \times .

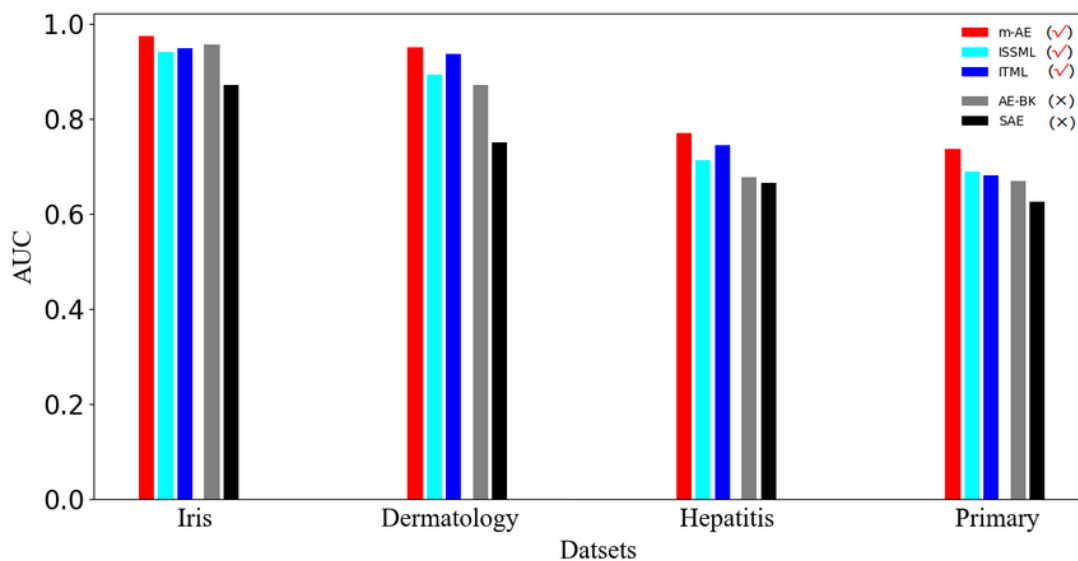


Figure 5

Results of ablation experiments

Comparisons between with distance metrics and with feature selection. The models with distance metrics are marked as the symbol \checkmark . The models with feature selection are marked as the symbol \neq .

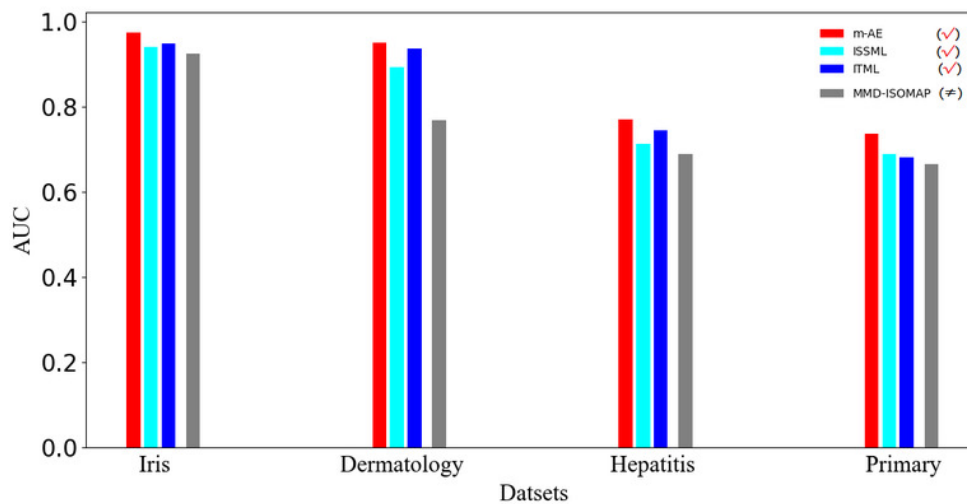


Figure 6

Visualization distributions.

The different extracted features are marked with different shapes and colors. The models with distance metrics are marked as the symbol \checkmark . The models with feature selection are marked as the symbol \neq . The models without both distance metrics and feature selection are marked as the symbol \times .

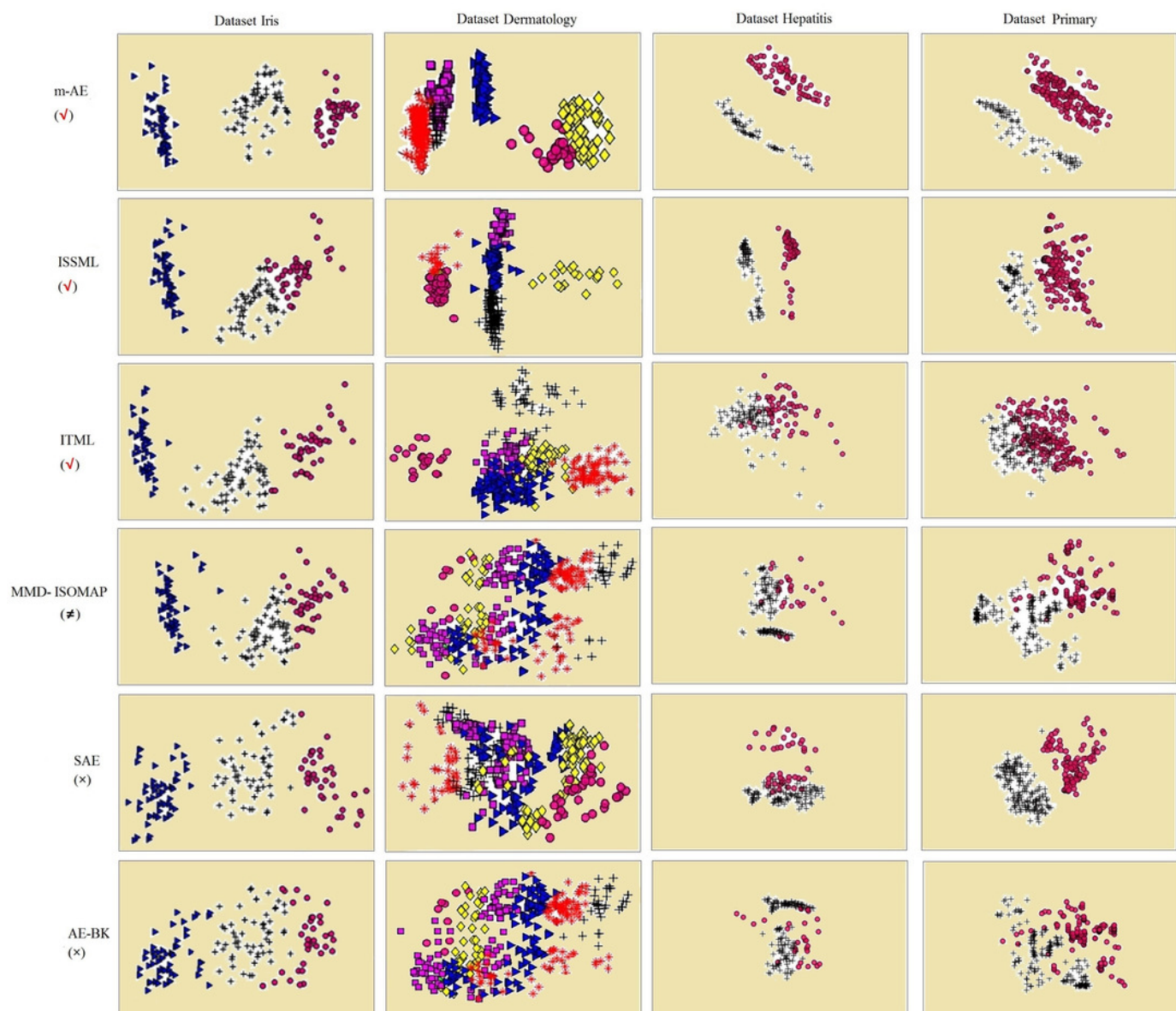
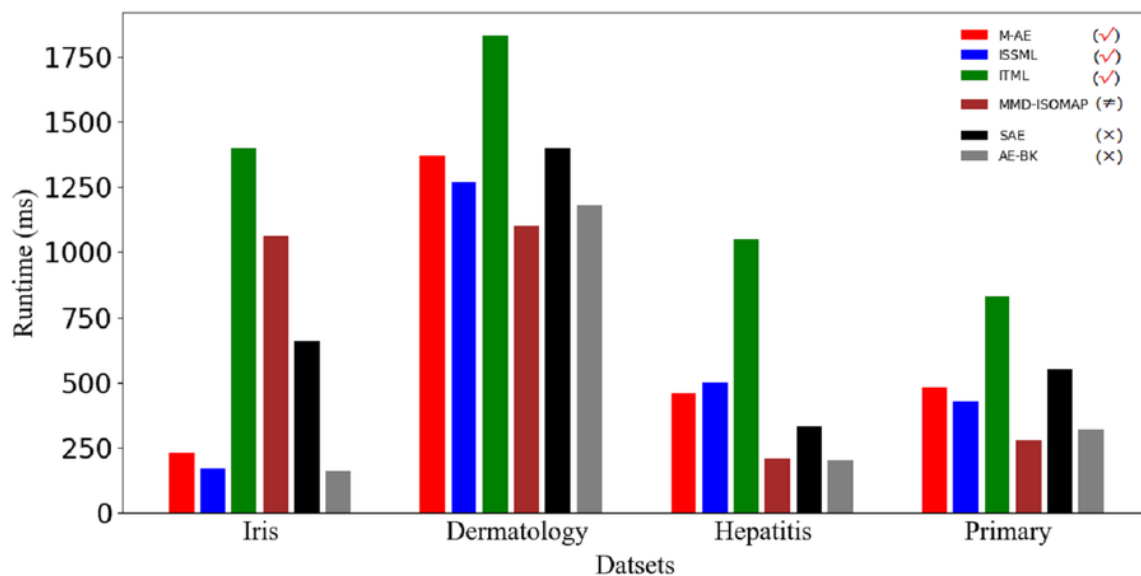


Figure 7

Runtime of methods.

The models with distance metrics are marked as the symbol \checkmark . The models with feature selection are marked as the symbol \neq . The models without both distance metrics and feature selection are marked as the symbol \times .



The models with distance metrics are marked as the symbol \checkmark . The models with feature selection are marked as the symbol \neq . The models without both distance metrics and feature selection are marked as the symbol \times .

Table 1(on next page)

Benchmark datasets

1

Dataset	Data volume	Data dimensionality	features
Iris	150	4	3
Primary	339	17	2
Hepatitis	155	19	2
Dermatology	366	33	6

2

Table 2 (on next page)

Accuracy of feature extraction.

The best accuracy for each dataset is shown in bold. The models using a distance metric are marked as the symbol \checkmark . The models using feature selection are marked as the symbol \neq .

The models without both a distance metric and feature selection are marked as the symbol \times .

1

		Iris	Dermatology	Hepatitis	Primary
m-AE	(✓)	0.9744 ± 0.0157	0.9506 ± 0.0137	0.7703 ± 0.0753	0.7375 ± 0.0534
ISSML	(✓)	0.9402 ± 0.0154	0.8931 ± 0.0284	0.7131 ± 0.0642	0.6886 ± 0.0865
ITML	(✓)	0.9488 ± 0.0120	0.9374 ± 0.0246	0.7457 ± 0.0622	0.6816 ± 0.0745
MMD-ISOMAP (≠)		0.9247 ± 0.0053	0.7680 ± 0.0377	0.6897 ± 0.0657	0.6664 ± 0.0733
SAE	(×)	0.9571 ± 0.0227	0.8707 ± 0.0892	0.6773 ± 0.0373	0.6700 ± 0.0166
AE-BK	(×)	0.8715 ± 0.1533	0.7511 ± 0.0099	0.6666 ± 0.0771	0.6252 ± 0.1052

2