

Identification of Autism spectrum disorder based on a novel feature selection method and Variational Autoencoder

Fangyu Zhang^a, Yanjie Wei^b, Jin Liu^d, Yanlin Wang^b, Wenhui Xi^b and Yi Pan ^{*b,c}

^aCollege of Engineering, Southern University of Science and Technology, Shenzhen, 518055, China

^bCentre for High Performance Computing, Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, 518055, China

^cCollege of Computer Science and Control Engineering, Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, 518055, China

^dSchool of Information Science and Engineering, Central South University, Changsha, 410083, China

Abstract

The development of noninvasive brain imaging such as resting-state functional magnetic resonance imaging (rs-fMRI) and its combination with AI algorithm provides a promising solution for the early diagnosis of Autism spectrum disorder (ASD). However, the performance of the current ASD classification based on rs-fMRI still needs to be improved. This paper introduces a classification framework to aid ASD diagnosis based on rs-fMRI. In the framework, we proposed a novel filter feature selection method based on the difference between step distribution curves (DSDC) to select remarkable functional connectivities (FCs) and utilized a multilayer perceptron (MLP) which was pretrained by a simplified Variational Autoencoder (VAE) for classification. We also designed a pipeline consisting of a normalization procedure and a modified hyperbolic tangent (tanh) activation function to replace the original tanh function, further improving the model accuracy. Our model was evaluated by 10 times 10-fold cross-validation and achieved an average accuracy of 78.12%, outperforming the state-of-the-art methods reported on the same dataset. Given the importance of sensitivity and specificity in disease diagnosis, two constraints were designed in our model which can improve the model's sensitivity and specificity by up to 9.32% and 10.21%, respectively. The added constraints allow our model to handle different application scenarios and can be used broadly.

Keywords: ASD; fMRI; Filter feature selection; VAE; ABIDE; Classification

1 Introduction

Autism spectrum disorder (ASD) is a common complex neurodevelopmental disorder that occurs in early childhood and has received much attention in recent years because of its high incidence and difficulty in curing. Although the manifestation of individuals with ASD varies greatly with age and ability [1], the disorder is characterized by core features in two areas—social communication and restricted, repetitive sensory-motor behaviors [2]. Some ASD patients can go undetected during childhood and be observed to have psychiatric comorbidity during adolescence [3], which presents challenges for traditional symptom-based diagnostic methods in ASD early diagnosis. Studies have shown that multiple biological factors can lead to the same ASD-related behavioral phenotype [4]. However, traditional symptom-based ASD

detection methods are unable to give a reliable diagnosis from a pathogenic perspective. To bridge this gap, non-invasive brain imaging techniques such as resting-state functional Magnetic Resonance Imaging (rs-fMRI) have been used to reveal valuable information about brain network organization [5] and contribute to a better understanding of the neural circuitry underlying ASD and its associated symptoms. By measuring changes in Blood Oxygen level-dependent (BOLD) signals, rs-fMRI can detect the functional connectivity patterns between the brain regions of interest (ROIs), and several studies have found abnormal functional connectivities in ASD subjects [6–8]. With the development of artificial intelligence technology, deep learning techniques have made it possible to process and analyze large amounts of fMRI data to discover patterns amongst the complex functional connectivities that are not apparent to the human eye and have achieved good results

*Corresponding author: Yi Pan (yi.pan@siat.ac.cn)

for ASD prediction, making fMRI-based deep learning models a promising auxiliary tool for ASD early screening.

However, small sample sizes and high feature dimensionality increase the challenge of developing a robust and well-performed machine learning model. One subject’s feature vector composed of functional connectivities (FCs) extracted from rs-fMRI usually has tens of thousands of dimensions which makes the model training a time-consuming process and contains a lot of noise resulting from the recording image process [9] and redundancy information which will affect the model performance. Feature selection is a dimensionality reduction technique that identifies the key features of a given problem [10]. The main goal of feature selection is to construct a subset of features as small as possible, which represent the vital features of the entire input data [11]. As an effective method to reduce training time and improve model performance, feature selection has been widely used in previous studies [12–17]. In the work of Guo et al. [12], a feature selection method based on multiple sparse autoencoders achieved a 9.09% model accuracy improvement based on 55 ASD and 55 Healthy Controls (HC) subjects from UM site of Autism Brain Imaging Data Exchange (ABIDE) I dataset [18]. Wang et al. [13] searched for informative features by SVM-RFE and obtained 90.60% accuracy by SVM classifier on a dataset consisting of 255 ASD and 276 HC subjects. In [14], a graph-based feature selection method was proposed to select remarkable FCs. Based on the refined features, a deep belief network (DBN) was trained and achieved a higher classification accuracy (76.4%) than previous studies on the entire ABIDE I dataset.

In previous studies, machine learning algorithms such as support vector machine (SVM), decision tree, and Gaussian naive Bayes have been applied to ASD recognition [19, 20], most of which belong to supervised learning methods. With the development of deep learning technology, pretraining classifiers using unsupervised deep learning methods has been proved to be helpful to improve the performance of classifiers. For instance, Heinsfeld et al. [21] utilized a denoising autoencoder to extract underlying representations of the input feature vectors and then trained fully connected layers for classifying ASD from HC based on 1035 subjects from the ABIDE I dataset. Their model achieved better classification accuracy (70%) than SVM and random forest (RF). Kong et al. [22] pretrained their model by a sparse autoencoder based on 182 subjects from NYU site of ABIDE I and achieved an accuracy of 90.39% and the

area under the receiver operating characteristic curve (AUC) of 0.9738 for ASD/HC classification. Although denoising autoencoders and sparse autoencoders are widely used in previous studies and got relatively good results, the setting of the noise ratios and sparsity parameters are subjective to some extent. Some other researchers such as Sherkatghanad et al. [23] and Shrivastava et al. [24] performed ASD classification by using CNN to process the functional connectivity matrix. These CNN-based methods can extract local characteristics of images, however, the functional connectivity matrix is not an image in Euclidean space and doesn’t have specific local characteristics. Other state-of-the-art methods that have been used in ASD classification include DBN [14], CapsNet [25], ASD-Diagnet [26], etc., all of these methods have achieved over 70% prediction accuracy. Given that the doctor’s diagnosis results can be affected by subjective factors such as different clinical experiences and fatigue, as well as objective factors such as the patient’s insignificant symptoms, false negative and false positive cases have always been difficult to avoid, while most previous studies primarily concentrated on developing high-accuracy models without taking measures to improve sensitivity or specificity, which are critical for reducing the false-negative rate and false-positive rate, respectively.

The purpose of this paper is to provide an ASD/HC classification framework based on rs-fMRI to improve ASD classification performance on heterogeneous datasets and to flexibly improve model sensitivity or specificity according to actual needs. In the framework, we proposed a novel feature selection method based on the difference between step distribution curves (DSDC) that not only contributed to higher classification accuracy but significantly speed up the training process. To get a more accurate and reliable ASD/HC classification model, we simplified the architecture of Variational Autoencoder (VAE) to pretrain the classifier and designed a pipeline consisting of a normalization and a modified hyperbolic tangent (tanh) activation function to replace the original tanh activation function, and adopted the threshold moving approach which is described in Section 2.4.3 to alleviate the impact of class imbalance of the dataset. In addition, we designed two constraints, by using which in the training process, model sensitivity or specificity can be effectively improved. The proposed method can potentially be used for ASD early screening and provide a valuable reference for doctors’ decision-making.

Our main contributions are summarized as follows:

1. We proposed an ASD/HC classification frame-

work including a novel feature selection method (the DSDC-based feature selection method), a simplified VAE pretraining method, and an MLP classifier. The accuracy of our classifier outperforms the state-of-the-art results reported on the same dataset with an outstanding training speed.

2. We designed two constraints that can be used during the model training process to effectively improve model sensitivity and specificity, respectively.

This paper is structured as follows: Section 2 introduces the dataset we used and rs-fMRI data preprocessing process (2.1-2.2), feature selection method (2.3) and details of our model (2.4-2.5). Section 3 discusses the experimental results and limitations. Finally, the conclusion and future work are presented in section 4.

2 Materials and methods

2.1 Participants

ABIDE I dataset is one of the most commonly used public datasets taken from 17 international sites (<http://preprocessed-connectomes-project.org/abide/>). In order to train a robust model with stronger generalization ability for the data from different sites, our study was carried out using all valid rs-fMRI data from ABIDE I including 505 ASD and 530 HC samples, the largest rs-fMRI subset of ABIDE I that has ever been used, the phenotype of which is summarized in Table 1.

Table 1: Demographic description of participants for ABIDE I

Site	ASD	HC	Male	Female	Subtotal	Average age
CALTECH	19	18	29	8	37	27
CMU	14	13	21	6	27	26
KKI	20	28	36	12	48	10
LEUVEN	29	34	55	8	63	18
MaxMun	24	28	48	4	52	25
NYU	75	100	139	36	175	15
OHSU	12	14	26	0	26	10
OLIN	19	15	29	5	34	16
PITT	29	27	48	8	56	18
SBL	15	15	30	0	30	34
SDSU	14	22	29	7	36	14
Stanford	19	20	31	8	39	9
Trinity	22	25	47	0	47	16
UCLA	54	44	86	12	98	13
UM	66	74	113	27	140	14
USM	46	25	71	0	71	22
Yale	28	28	40	16	56	12
Total:	505	530	878	157	1035	

Four different preprocessed datasets have been provided by ABIDE I according to four pipelines (CPAC [27], CCS, DPARSE, and NIAK). The CPAC pipeline was considered in our work, partly because previous studies such as Zhang et al. [28] have compared the four different pipelines and found that data preprocessed with CPAC pipeline can achieve better classification results, and partly because it allows our model to be compared and evaluated against most of the other methods that have chosen CPAC pipeline as well. In addition, ABIDE I dataset provides data preprocessed by seven brain atlases among which Craddock 200 (CC200) [29] and Craddock 400 (CC400) have been proved by previous studies better than other atlases such as Automated Anatomical Labeling(AAL), Dosenbach160, etc [26, 30]. CC400 atlas was adopted in our work mainly because CC400 has a more detailed division of ROIs than CC200. Another important reason is that the DSDC-based feature selection greatly reduced the dimension of the original feature vector, which enables us to complete the model training in a short time even based on a complex brain atlas like CC400 (392 ROIs, 76636 features for each subject).

2.2 Functional Connectivity Measures and Subject’s Feature Vector

CC400 is a brain atlas with 392 ROIs from which 392 time series were extracted. The Pearson correlation coefficient (PCC) was calculated between each pair of time series by Formula (1) to measure the coactivation level between each pair of ROIs and form a PCC matrix. For each subject, we took the upper triangle of the PCC matrix and removed the main diagonal elements, after which the remaining triangle was flattened into a one-dimensional feature vector including $392 \times (392 - 1) / 2 = 76636$ features. The entire process of generating subjects’ feature vectors is shown in Figure 1.

$$PCC_{x,y} = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}} \quad (1)$$

Where $PCC_{x,y}$ is the PCC between time series x and y ; N is the length of time series; \bar{x} and \bar{y} are the mean value of time series x and y .

2.3 DSDC-based feature selection

A subject can be represented by a feature vector as explained in section 2.2, while most features have complex distributions across subjects. In order to reduce

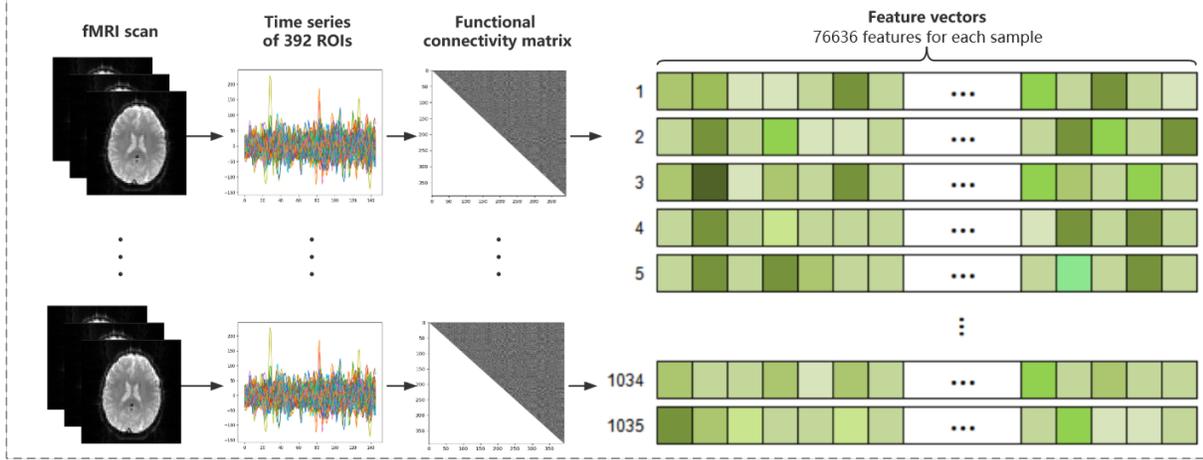


Figure 1: The process of generating subjects' feature vectors from fMRI images

the complexity of the distribution curves, the value range of features was evenly divided into 20 subintervals. Within each subinterval, for each class, the number of samples was divided by the total sample size of the class and got a normalized value. By using all the normalized values, step distribution curves (see Figure 2 (B)) were created to approximate the original feature distribution curves (see Figure 2 (A)). We defined the DSDC score in equation (2) to measure the distribution difference of positive and negative samples.

$$DSDC_score = \sum_{i=b_0+\delta}^{b_1} |n_i^+ / N^+ - n_i^- / N^-| \quad (2)$$

Where b_0 and b_1 represent the lower bound and upper bound of feature values; δ is the span of subinterval; n_i^+ and n_i^- are the numbers of positive and negative samples whose feature values are in $[i - \delta, i)$; N^+ and N^- are the positive and negative sample sizes. The larger the DSDC score of a feature, the more discriminative the feature is.

In our experiment, a feature was considered remarkable if its DSDC score is bigger than a preset filter threshold (0.241). The filter threshold was determined by the following process: first, 55 feature subsets of different sizes were generated by the DSDC-based feature selection method with different filter thresholds, then an MLP with two hidden layers was used to perform 10-fold cross-validation on each feature subset and the average accuracy was calculated. We chose the filter threshold corresponding to the highest average accuracy as the preset filter threshold (see Figure 3). Through DSDC-based feature selection,

3170 remarkable features were selected from the original 76636 features whose dimension was reduced by 95.86%. Subsequent experimental results show that the feature selection not only improved the classifier's accuracy but also greatly reduced the training time of the deep learning model. In addition, the feature selection process of ABIDE I dataset with the input matrix size of 1035×76636 takes 28.12 seconds based on Intel Xeon Silver 4114 CPU, which reflects that the DSDC score is computationally efficient.

2.4 Deep neural network classifier

This section introduces the architecture of our classifier and its training process.

2.4.1 Simplified VAE and MLP

The training process of our classifier consists of simplified VAE unsupervised pretraining and MLP supervised fine-tuning.

VAE [31, 32] is a generative model with an encoder and a decoder. Given an input vector x , the VAE's encoder outputs the latent space parameters (μ and $\log var$) by equations (3), (4) and (5). Then the parameters are used to generate the latent variables z by equation (6). At last, VAE's decoder uses the latent variables to reconstruct the input vector.

$$h = W_1 x + b_1 \quad (3)$$

$$\mu = W_2 h + b_2 \quad (4)$$

$$\log var = W_3 h + b_3 \quad (5)$$

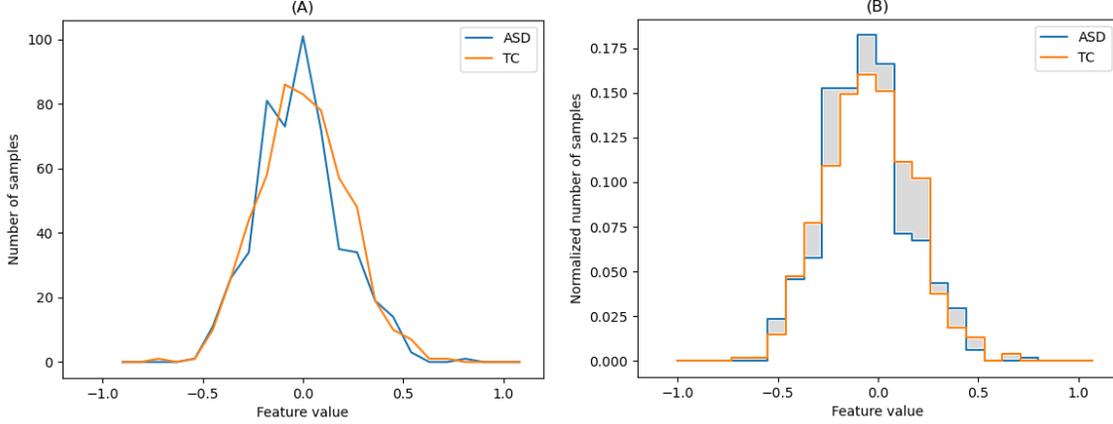


Figure 2: The original distribution curves and the corresponding step distribution curves for the functional connectivities between Right Middle Frontal Gyrus and Right Superior Frontal Gyrus. (A) Original distribution curves; (B) Step distribution curves

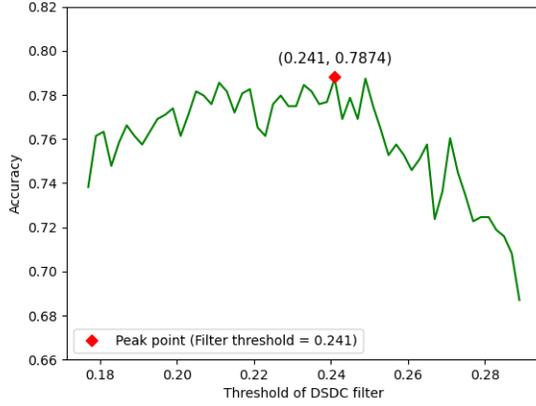


Figure 3: Selecting the filter threshold corresponding to the highest average accuracy of the 10-fold cross-validation

$$z = \mu + \epsilon \times e^{0.5 \times \log var} \quad (6)$$

where h is the output of the first hidden layer of VAE's encoder, μ and $\log var$ are parameters of the latent space, $\{W_1, W_2, W_3, b_1, b_2, b_3\}$ are weights and biases of the encoder, ϵ is a random number sampled from $N(0, 1)$ distribution.

We use the root-mean-square propagation (RMSProp) [33] as the backpropagation method to optimize the loss function of VAE which is defined as follows,

$$Loss(W, b) = MAE + \beta \sum_{i=1}^n KL[N(\mu_i, e^{\log var_i}) \parallel N(0, 1)] \quad (7)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |f_{W,b}(x_i) - x_i| \quad (8)$$

$$KL[N(\mu_i, e^{\log var_i}) \parallel N(0, 1)] = -\frac{1}{2}(1 + \log var_i - e^{\log var_i} - \mu_i^2) \quad (9)$$

Where $Loss(W, b)$ represents the loss function of VAE. MAE represents the mean absolute error [34] between prediction and true label, n is the number of samples, x_i is the input of the i th sample, $f_{W,b}(x_i)$ is output of VAE's decoder. $KL[N(\mu_i, e^{\log var_i}) \parallel N(0, 1)]$ represents the Kullback-Leibler divergence between $N(\mu_i, e^{\log var_i})$ and $N(0, 1)$, μ_i and $\log var_i$ are parameters of latent space generated by encoder of VAE.

As seen from equations (4) and (5), the μ and $\log var$ are generated by two branch networks with different parameters, respectively. However, our purpose is transferring the parameters of the pretrained VAE's encoder to MLP for fine-tuning, thus we simplified the structure of VAE's encoder by using a unified network to generate μ and $\log var$ simultaneously following equation (10) to ensure the encoder's structure is the same as the MLP's structure.

$$\mu = \log var = W_2 h + b_2 \quad (10)$$

After pretraining, we use the parameters of the VAE's encoder as the initialization parameters of the MLP and fine-tune the MLP's parameters through supervised training. Cross entropy is used as the loss function and RMSProp is adopted as the backpropagation method. A softmax layer is added after the MLP to calculate probabilities to determine the label of each subject. Afterwards, evaluation metrics such

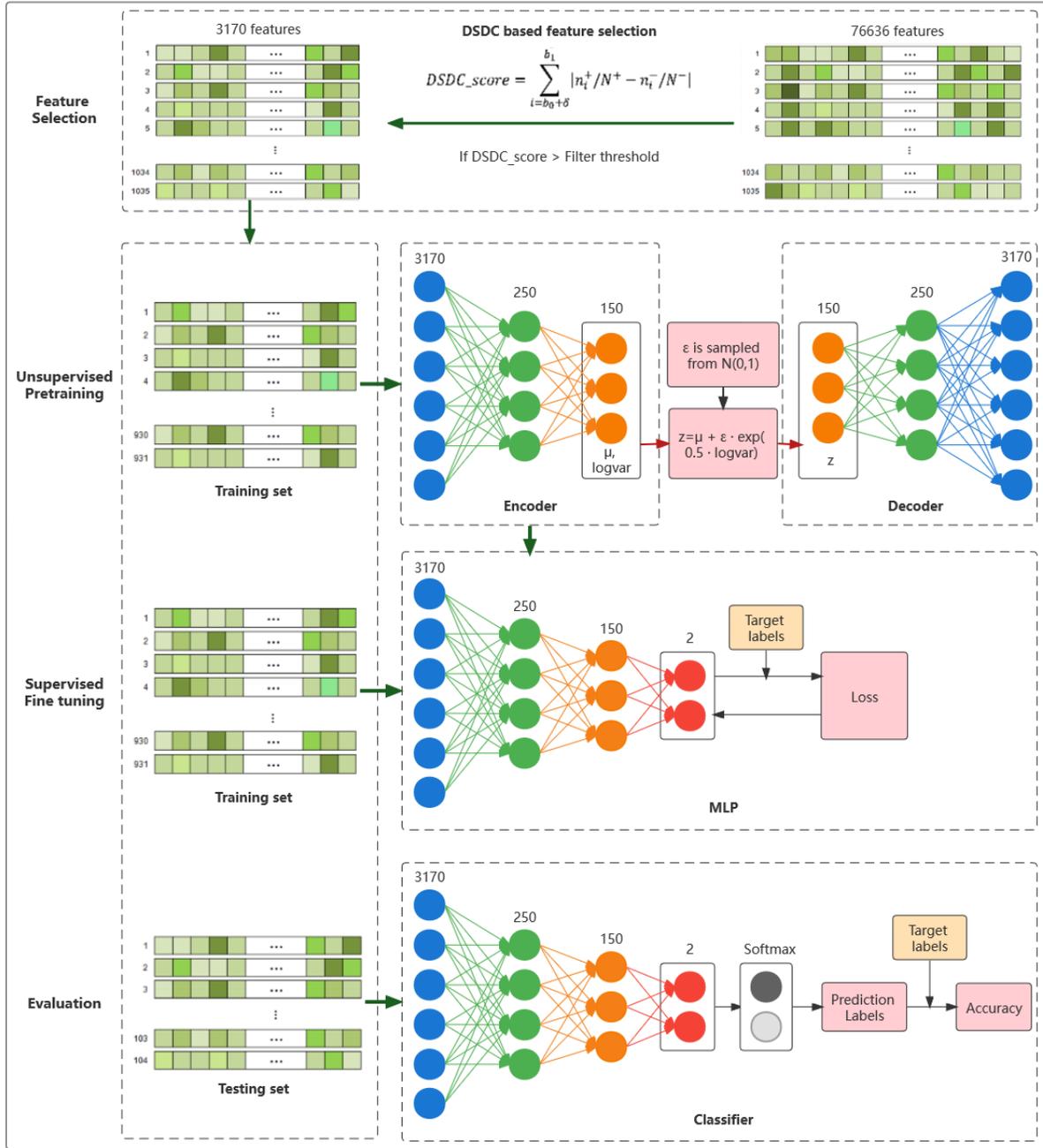


Figure 4: The entire process of feature selection, pretraining, fine-tuning and model evaluation

as accuracy, sensitivity, and specificity can be calculated by using the prediction label and true label. The model’s architecture and the entire process of feature selection, pretraining, fine-tuning, and evaluation are shown in Figure 4.

2.4.2 Normalization and modified tanh activation function

To further improve the classification accuracy, a pipeline consisting of a normalization procedure and a modified tanh activation function was designed to replace the original tanh activation function.

The normalization can be performed according to Equation (11) through which the outputs of hidden layers are mapped to $[-1, 1]$ in order to match the

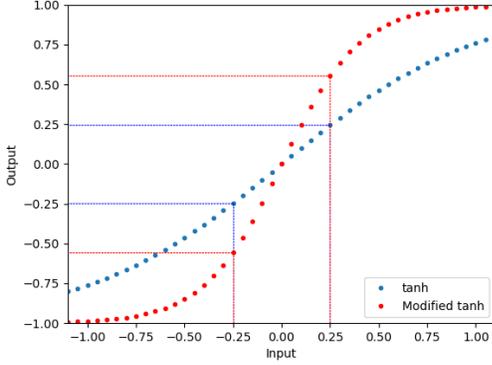


Figure 5: A comparison of the tanh and the modified tanh

subsequent activation function.

$$x_{norm} = \frac{2(x - x_{min})}{x_{max} - x_{min}} - 1 \quad (11)$$

Where x is the output of one hidden node and x_{norm} is the normalized output; x_{max} and x_{min} is the maximum and minimum output of the hidden node.

After normalization, the modified tanh activation function in Equation (12) is applied to the outputs of hidden layers. Compared with tanh, the modified tanh can make the features whose values close to zero more discriminating by mapping them to a larger interval (see Figure 5)

$$y_{norm_act} = \frac{e^{2.5 \cdot x_{norm}} - e^{-2.5 \cdot x_{norm}}}{e^{2.5 \cdot x_{norm}} + e^{-2.5 \cdot x_{norm}}} \quad (12)$$

Where y_{norm_act} is the output of normalization-activation pipeline.

2.4.3 Threshold moving

Threshold moving is an approach to alleviate the impact of class imbalance on classification performance by adjusting the classification threshold. When the data is balanced, the classification threshold of a softmax layer's output is usually set to 0.5. However, the number of HC samples is larger than that of ASD samples in ABIDE I, therefore, the threshold moving approach was adopted in our work and the decision is made according to the following rules.

$$\frac{p_{ASD}}{p_{HC}} \cdot \frac{N_{HC}}{N_{ASD}} > 1 \Leftrightarrow \frac{p_{ASD}}{p_{HC}} > \frac{N_{ASD}}{N_{HC}} \Rightarrow ASD \quad (13)$$

$$\frac{p_{ASD}}{p_{HC}} \cdot \frac{N_{HC}}{N_{ASD}} < 1 \Leftrightarrow \frac{p_{ASD}}{p_{HC}} < \frac{N_{ASD}}{N_{HC}} \Rightarrow HC \quad (14)$$

Where $[p_{ASD}, p_{HC}]$ are the outputs of the softmax layer, representing the probabilities of a sample being ASD and HC. N_{HC} and N_{ASD} are the number of HC and ASD samples in the training set, respectively.

As shown in equations (13) and (14), in the classification decision stage, we increase the weight of ASD which is equivalent to relaxing the criterion for determining ASD.

2.5 Additional constraints during model training

This section will introduce two constraints that can help to train models with higher sensitivity and specificity, respectively. The constraints are used in the model training process to determine whether the optimized model parameters of a certain training epoch should be saved. The mechanism of the constraints is detailed in Algorithms 1.

Take constraint 1 for example, in the training process, constraint 1 helps to improve the difference between sensitivity and specificity while improving the accuracy. The sensitivity increases because it is positively correlated with the difference between sensitivity and specificity which can be proved as follows.

Because:

$$Accuracy = \frac{TP + TN}{N_{ASD} + N_{HC}} \quad (15)$$

$$Sensitivity = \frac{TP}{N_{ASD}}, \quad Specificity = \frac{TN}{N_{HC}} \quad (16)$$

$$Sensitivity - Specificity = \frac{TP}{N_{ASD}} - \frac{TN}{N_{HC}} \quad (17)$$

$$N_{ASD}, N_{HC} = Constant \quad (18)$$

Assume:

$$Accuracy = Constant \quad (19)$$

Then:

$$TP + TN = Constant \quad (20)$$

Obviously:

$$\frac{TP}{N_{ASD}} - \frac{TN}{N_{HC}} \uparrow \Rightarrow TP \uparrow \Rightarrow Sensitivity \uparrow \quad (21)$$

But if the difference between sensitivity and specificity keeps increasing and exceeds a certain value, the model's performance will deteriorate. Therefore an appropriate threshold is needed and the selection of the threshold will be discussed in detail in Section 3.3.

3 RESULTS AND DISCUSSION

On ABIDE I dataset (505 ASD / 530 HC), most previous studies evaluated their models through a single 10-fold cross-validation, the evaluation results of which may be susceptible to the randomness of the dataset division. In our work, 10-fold cross-validation was repeated 10 times to evaluate our model more objectively. For each 10-fold cross-validation, we performed a stratified sampling to randomly divide the dataset into a training set, a validation set, and a testing set in an 8:1:1 ratio. The training set was used to train a model and the validation set was used to stop training automatically when the validation accuracy stopped increasing to avoid overfitting, and the testing set is used for evaluating the classification performance of the trained model. Accuracy, sensitivity, specificity, and training time are obtained by calculating the mean value of 10 times experiments for model evaluation. A grid search with a step size of 50 was performed to optimize the model’s layer configuration from full[100]-full[100]-full[2] to full[1000]-full[900]-full[2], through which the layer configuration is determined as full[250]-full[150]-full[2], corresponding to the highest accuracy (78.12%), where full[N] denotes a fully-connected layer with N outputs.

In the following subsections, first, the DSDC-based feature selection method is compared with two other commonly used feature selection methods to highlight the advantage of the former. Next, the contribution of each procedure in our framework to classification accuracy is discussed. Then, we analyze the model performance by using constraints with different thresholds and provide a feasible threshold selection range. After that, our experimental results are compared with other state-of-the-art studies on the same dataset. Finally, the limitations of the current work are discussed.

3.1 Evaluation of feature selection methods

The advantage of the proposed DSDC-based feature selection method is highlighted by comparing it with two widely used feature selection methods based on F-score [35] and PCC. At first, all features are ranked by DSDC score, F-score, and the absolute value of PCC in descending order. Then top n (n varies from 500 to 30000) features of the different feature rankings were fed into the same SVM classifier. Finally, we compared the accuracy of the SVM classifier under different feature rankings (Figure 6).

As in Figure 6, the red curve (representing DSDC) rises faster, indicating the top features of DSDC-based

Algorithm 1

Input: 1.*initial model* \triangleright Untrained Model with random initialization of parameters
2.*training set* \triangleright The dataset used to train a model
3.*validation set* \triangleright The dataset used to stop training automatically
4.*constraint_type* \triangleright Parameter used to choose Constraint 1 or Constraint 2
5.*threshold* \triangleright A preset threshold for constraints, the effect of which will be discussed in Section 3.4

Output: *final model* \triangleright The trained model

```

1:  $max\_acc = 0$ 
2:  $\delta = -1$ 
3: for  $training\_epoch = 1$  to  $max\_training\_epoch$  do
4:   Training model on training set  $\rightarrow$  current model
5:   Calculating accuracy, sensitivity and specificity of current model on validation set  $\rightarrow v\_acc, v\_sen, v\_spe$ 
6:   if  $constraint\_type == 1$  then  $\triangleright$  Constraint 1
7:     if  $v\_acc \geq max\_acc$  and  $v\_sen - v\_spe \geq \delta$  then
8:        $max\_acc = v\_acc$ 
9:       if  $\delta < threshold$  then
10:         $\delta = v\_sen - v\_spe$ 
11:      else
12:         $\delta = threshold$ 
13:       $final\_model = current\_model$ 
14:   else if  $constraint\_type == 2$  then  $\triangleright$  Constraint 2
15:     if  $v\_acc \geq max\_acc$  and  $v\_spe - v\_sen \geq \delta$  then
16:        $max\_acc = v\_acc$ 
17:       if  $\delta < threshold$  then
18:         $\delta = v\_spe - v\_sen$ 
19:      else
20:         $\delta = threshold$ 
21:       $final\_model = current\_model$ 
22:   else
23:     if  $v\_acc \geq max\_acc$  and  $|v\_sen - v\_spe| \leq \delta$  then
24:        $max\_acc = v\_acc$ 
25:       if  $\delta > threshold$  then
26:         $\delta = |v\_sen - v\_spe|$ 
27:      else
28:         $\delta = threshold$ 
29:       $final\_model = current\_model$ 
30: return  $final\_model$ 

```

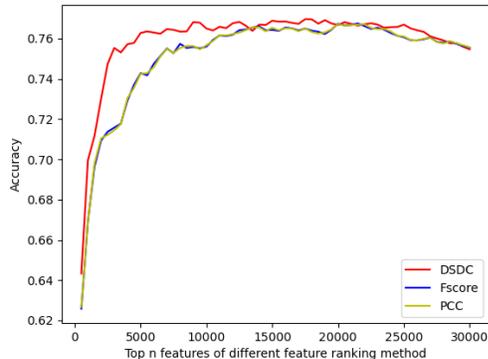


Figure 6: A comparison of different feature selection methods based on the same SVM classifier

feature selection include more vital information than the top features of the other two methods (F-score and PCC). With the increase in the number of top features, more important features have been selected by all the three methods so that the accuracy difference between the three methods becomes smaller, while the DSDC-based feature selection method still showed better performance. This illustrates that DSDC can help to select more refined features and contribute to better classification performance. It is worth mentioning that we use an SVM classifier to compare different feature selection methods instead of a deep learning model because the model parameters will change according to different input dimensions. Changes in the parameters of a deep learning model commonly have an influence on the model performance, and it is difficult to tell whether changes in classification performance are mainly caused by the new inputs or the changes in model parameters.

3.2 The performance analysis of our framework

The contribution to the classification accuracy of each procedure in our framework is shown in Figure 7. Besides, we performed two-sample t-tests and calculated p-values to investigate the significance of accuracy improvement (p-value < 0.05 indicates a significant improvement in accuracy). Due to the feature vectors without dimensionality reduction can bring a huge time cost to the training and hyperparameter tuning of a deep learning model, an SVM classifier is used instead to evaluate the performance improvement contributed by DSDC-based feature selection (see the 1st and 2nd rows). It can be observed that the DSDC-based feature selection resulted in a significant accuracy improvement of 5.56% (p-value < 0.01). The 2nd and 3rd rows verify the advantages

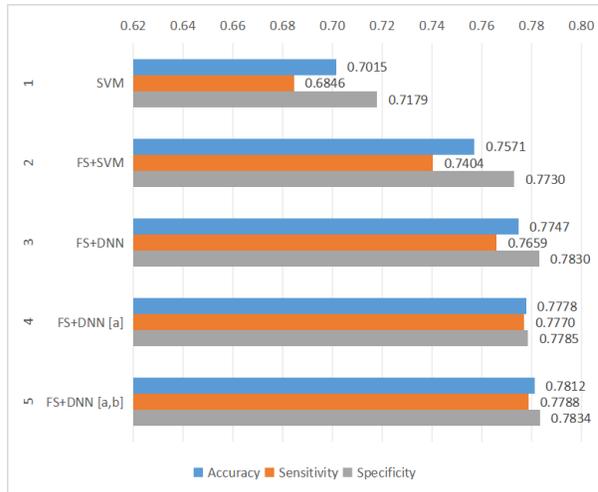


Figure 7: Performance comparison of SVM, FS+SVM, FS+DNN, FS+DNN[a] and FS+DNN[a,b] models

- FS: The DSDC-based feature selection described in Section 2.3
- DNN: MLP pretrained by simplified VAE described in Section 2.4.1
- a: Using normalization and modified tanh activation function described in Section 2.4.2 instead of tanh in DNN
- b: Using threshold moving approach described in Section 2.4.3 in DNN

of deep learning over traditional machine learning by comparing the deep learning model pretrained by simplified VAE with SVM (p-value < 0.01). By replacing the original tanh activation function with the normalization-activation pipeline described in section 2.4.2 and using the threshold moving approach described in section 2.4.3, a 0.65% increase in accuracy is observed by comparing the 3rd and 5th rows (p-value = 0.02), indicating that using these methods can significantly improve the accuracy of our model, while the individual effects of the normalization-activation pipeline (0.31% accuracy increase showed in 3th and 4th rows) and threshold moving (0.34% accuracy increase showed in 4th and 5th rows) are slight. In the loss function of the simplified VAE, we have also tried to use MSE loss instead of MAE loss, whereas the accuracy of the final classifier slightly decreased by 0.28%. One possible reason is that the MAE loss is more robust to outliers than the MSE loss, and more than 89% of the 3170 input features have outliers if values more than three times the standard deviation away from the mean value is regarded as outliers.

Figure 8 was used to investigate the influence of the simplified VAE pretraining method on the convergence speed of the classifier by comparing the average training accuracy of the 10 times 10 folds' experiments after each training epoch. As shown in Figure 8, the red curve (representing the MLP pretrained by simplified VAE) is above the blue curve (representing the

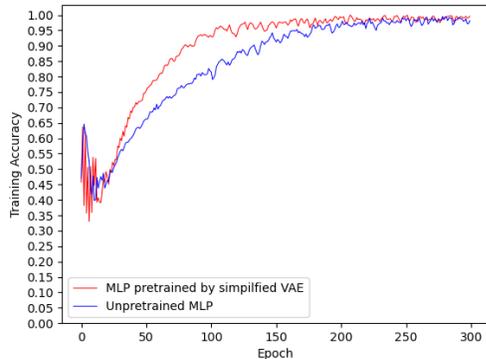


Figure 8: A comparison of convergence speed between MLP pretrained by simplified VAE and unpretrained MLP

unpretrained MLP), which shows that the simplified VAE pretraining can speed up the convergence of the MLP by providing it with initialization parameters.

Since different sites of ABIDE I adopt different scanning protocols, there is heterogeneity among data from different sites. To investigate whether the simplified VAE pretraining method can extract more useful information from a large heterogeneous dataset. We selected the five largest sites from the ABIDE I dataset and performed pretraining, fine-tuning, and inference independently on each of them (the corresponding accuracy is represented by blue bars in Figure 9). For comparison, we used the entire ABIDE I dataset (excluding the testing set) for pretraining and performed fine-tuning and inference independently on each of the 5 sites (the corresponding accuracy is represented by orange bars in Figure 9). As shown in Figure 9, for LEUVEN, UCLA, UM, and NYU sites, classifiers pretrained on the entire ABIDE I dataset achieved higher accuracy than those pretrained on a single site, which indicates that in most cases more useful patterns could be extracted from a larger heterogeneous dataset through the simplified VAE. However, the result of the USM site is an exception. One possible reason is the features' distribution the USM samples has a relatively bigger difference from that of other sites.

3.3 The effects of constraints on model performance

Section 2.5 has described the mechanism of the two constraints used to train models with higher sensitivity and specificity and explained that an appropriate threshold is necessary to ensure a good model performance. This section will analyze the influence of constraints' thresholds on model performance. Our

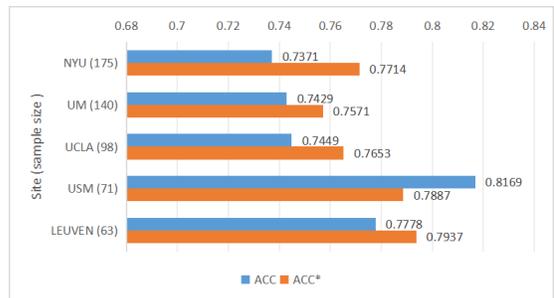


Figure 9: Accuracy comparison of MLPs pretrained independently on the five largest sites of ABIDE I dataset and pretrained on the entire ABIDE I dataset

- ACC: The accuracy of MLP pretrained and fine-tuned on a single site
- ACC*: The accuracy of MLP pretrained on the entire ABIDE I dataset and fine-tuned on a single site

experimental results demonstrate that the constraints can significantly improve the model's sensitivity or specificity at a cost of small accuracy reduction while ensuring the overall performance. Figure 10 summarizes the trade-off among accuracy, sensitivity, and specificity under different thresholds in Algorithm 1. Take constraint 1 as an example, in the range of 0 to 0.3, selecting a higher threshold will improve the sensitivity, but at a cost, the accuracy and specificity will decrease; When the threshold exceeds 0.3, the performance of the model deteriorates as the threshold increases. Constraint 2 follows a similar pattern. Therefore, the value range for the threshold is 0 to 0.3. In our work, we set the threshold to be 0.3 and trained two models with 87.20% sensitivity and 88.55% specificity by using constraint 1 and constraint 2 in Algorithm 1, respectively. Compared with the model trained without using the two constraints (Model 1 in Table 2), Model 2 in Table 2 improved sensitivity by 9.32% at the cost of an accuracy reduction of 2.76% and Model 3 in Table 2 improved specificity by 10.21% at the cost of an accuracy reduction of 4.38%. For disease diagnosis, a model with high sensitivity and specificity can reduce missed diagnosis and misdiagnosis rates which reflects our model could potentially adapt to the different cases for clinical application. For example, a model with higher sensitivity is very useful for the cases like COVID-19 screening since missing the diagnosis of a virulent communicable illness has far-reaching public health implications and may accelerate the pandemic. Models with higher sensitivity and specificity can also be used to perform a double check on the doctor's diagnosis results. For subjects whose doctor's diagnosis results are inconsistent with the model's prediction results, further analysis or detection can be performed to eliminate

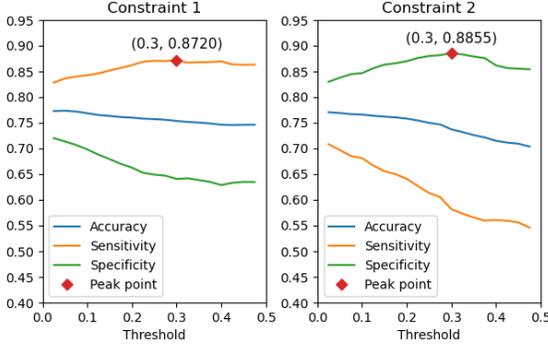


Figure 10: Influence of different thresholds for constraints in Algorithm 1 on accuracy, sensitivity and specificity of classifier

the potential diagnostic errors. In order to evaluate the influence of Algorithm 1 on the model’s overall performance, Figure 11 compares the models’ average AUCs, ROCs, and DET curves before and after using Algorithm 1. The result shows that the three models have similar ROCs and AUCs which demonstrates that the proposed constraints can ensure the model’s overall performance while improving sensitivity or specificity. In Figure 11, the DET curves of the classifier trained without using Algorithm 1 and classifiers trained with constraint 1 and constraint 2 in Algorithm 1 are represented by green, red, and blue curves, respectively. It can be observed that the red curve is generally below the green curve when the false-negative rate is lower than 20% and the blue curve is generally below the green curve when the false-positive rate is lower than 20%. This indicates that when the false-negative rate and the false-positive rate are relatively low, constraint 1 and constraint 2 can contribute to better model performance, respectively.

3.4 Comparison with other state-of-the-art methods on the same dataset

In this section, we make a comparison of our proposed method with several previous methods based on the same dataset in Table 2. The state-of-the-art methods used for comparison include denoising autoencoder, convolutional neural network (CNN), DBN, CapsNet, and ASD-Diagnet, all of which have achieved over 70% prediction accuracy. In our work, we trained three models based on simplified VAE and MLP, Model 2 and 3 were trained by using the constraint 1 and constraint 2 described in algorithm 1, respectively, while Model 1 is trained without using the two constraints. The accuracy (78.12%), sensitivity (87.20%), and specificity (88.55%) of Model 1, Model

2, and Model 3 exceed the corresponding results in previous studies on the same dataset, respectively. This highlights the outstanding classification performance and flexibility of our framework. Our method also runs efficiently due to less number of selected features. The last column of Table 2 lists the training time and the corresponding GPU used for training. The training time of our model is 85s for 10-fold cross-validation (8.5 seconds for training a single model), which has an advantage over other studies after taking into account the GPUs’ performance differences. A model with a fast training speed can be retrained in a short time as new subjects arrive, which means that the model can potentially achieve real-time optimization and improvement in practical use.

3.5 Limitations

In the present study, the experimental results are based on the ABIDE I dataset. More ASD datasets or other neurological diseases are expected to be used to evaluate the classification framework in the future. In addition, two other modalities of MRI (structural MRI and diffusion tensor imaging) have not been used in our work. Since different MRI modalities contain complementary information for ASD identification, fusing multiple modalities for ASD classification may work better than just using rs-fMRI. The proposed deep learning model can potentially be used to discover vital functional connectivities by studying its explainability, which can help to understand the ASD mechanism. However, this has not been included in the current work.

4 Conclusion

In this study, we proposed a novel filter feature selection method (DSDC-based feature selection method) to select remarkable FCs and designed a deep learning model with two procedures – simplified VAE pretraining and MLP fine-tuning. In our model, we designed a pipeline consisting of normalization and a modified tanh activation function to replace the original tanh function and adopted the threshold moving approach, which can further improve the classification accuracy. In addition, we proposed two constraints that can help to train models with higher sensitivity or specificity. The outstanding classification performance on the heterogeneous dataset and adjustable sensitivity and specificity suggest that our method goes one step further based on state-of-the-art methods and has the potential to be a viable auxiliary method for ASD early detection.

Table 2: Performance comparison between our work and other state-of-the-art studies

Model	Dataset	Model	Validation	ACC	SEN	SPE	Best cases in 10 times 10-fold CV			Worst cases in 10 times 10-fold CV			Training time of 10-fold CV
							ACC	SEN	SPE	ACC	SEN	SPE	
Model 1 (Our proposed)	ABIDE I 505 ASD & 530 HC	Simplified VAE + MLP	Mean value of 10 times 10-fold CV	0.7812	0.7788	0.7834	0.7894	0.804	0.8019	0.777	0.7564	0.7528	85s NVIDIA Tesla P100
Model 2 (Our proposed)	ABIDE I 505 ASD & 530 HC	Simplified VAE + MLP using constraint 1 in Algorithm 1	Mean value of 10 times 10-fold CV	0.7536	0.8720	0.6408	0.7739	0.895	0.6943	0.7391	0.8535	0.6038	85s NVIDIA Tesla P100
Model 3 (Our proposed)	ABIDE I 505 ASD & 530 HC	Simplified VAE + MLP using constraint 2 in Algorithm 1	Mean value of 10 times 10-fold CV	0.7374	0.5820	0.8855	0.7536	0.6198	0.9113	0.7043	0.5347	0.866	85s NVIDIA Tesla P100
Huang et al. [14]	ABIDE I 505 ASD & 530 HC	DBN with DE-based optimizer	10-fold CV	0.7640 ±0.022	–	–	–	–	–	–	–	–	–
Shrivastava et al. [24]	ABIDE I 505 ASD & 530 HC	CNN	10-fold CV	0.7602	0.7004	0.8169	–	–	–	–	–	–	–
Jiao et al. [25]	ABIDE I 505 ASD & 530 HC	CapsNet	10-fold CV	0.7100	0.7300	0.6600	–	–	–	–	–	–	–
Eslami et al. [26]	ABIDE I 505 ASD & 530 HC	ASD-DiagNet	10-fold CV	0.7030	0.6830	0.7220	–	–	–	–	–	–	41 min NVIDIA Tesla K40c
Shekarghanad et al. [23]	ABIDE I 505 ASD & 530 HC	CNN	10-fold CV	0.7022	0.7746	0.6182	–	–	–	–	–	–	12h and 30min NVIDIA Tesla K80
Heinsfeld et al. [21]	ABIDE I 505 ASD & 530 HC	SAE	10-fold CV	0.7000	0.7400	0.6300	–	–	–	–	–	–	32h 52 m 36 s NVIDIA Tesla K40

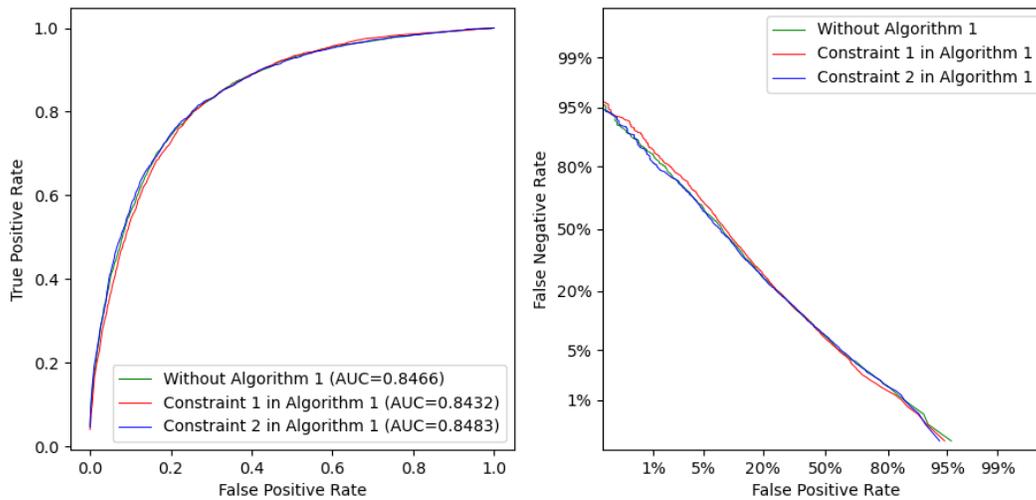


Figure 11: ROCs (left) and DET curves (right) of classifiers trained with and without constraints in Algorithm 1

- The threshold of the constraints is set to 0.3

Future work will focus on fusing multimodal MRI data to utilize the complementary information for ASD identification, which is expected to further improve the classification performance, and experiments will be performed on more datasets to verify the robustness of the model. In addition, further work will be carried out to study the explainability of the proposed model, which can potentially be used to discover vital functional connectivities and help to understand the ASD mechanism.

Acknowledgments

This work was partly supported by the Shenzhen KQTD Project No. KQTD20200820113106007, National Key Research and Development Program of China under Grant No. 2018YFB0204403, Strategic Priority CAS Project XDB38050100, National Science Foundation of China under grant no. U1813203, the Shenzhen Basic Research Fund under grant no. RCYX2020071411473419 and JSJG20201102163800001, CAS Key Lab under grant no. 2011DP173015.

References

[1] Uta Frith and Francesca Happé. Autism spectrum disorder. *Current biology*, 15(19), 2005.

[2] Catherine Lord, Mayada Elsabbagh, Gillian Baird, and Jeremy Veenstra-Vanderweele.

Autism spectrum disorder. *The Lancet*, 392(10146):508–520, 2018.

[3] Shilpa Aggarwal and Beth Angus. Misdiagnosis versus missed diagnosis: diagnosing autism spectrum disorder in adolescents. *Australasian Psychiatry*, 23(2):120–123, 2015.

[4] Daniel H Geschwind and Pat Levitt. Autism spectrum disorders: developmental disconnection syndromes. *Current Opinion in Neurobiology*, 17(1):103–111, 2007. Development.

[5] Sergey M. Plis, Michael P. Weisend, Eswar Damaraju, Tom Eichele, Andy Mayer, Vincent P. Clark, Terran Lane, and Vince D. Calhoun. Effective connectivity analysis of fmri and meg data collected under identical paradigms. *Computers in Biology and Medicine*, 41(12):1156–1165, 2011. Special Issue on Techniques for Measuring Brain Connectivity.

[6] Jocelyn V Hull, Lisa B Dokovna, Zachary J Jaccokes, Carinna M Torgerson, Andrei Irimia, and John Darrell Van Horn. Resting-state functional connectivity in autism spectrum disorders: a review. *Frontiers in psychiatry*, 7:205, 2017.

[7] Kaustubh Supekar, Lucina Q Uddin, Amirah Khouzam, Jennifer Phillips, William D Gaillard, Lauren E Kenworthy, Benjamin E Yerys, Chandan J Vaidya, and Vinod Menon. Brain hypercon-

- nectivity in children with autism and its links to social deficits. *Cell reports*, 5(3):738–747, 2013.
- [8] Marcel Adam Just, Timothy A Keller, Vicente L Malave, Rajesh K Kana, and Sashank Varma. Autism as a neural systems disorder: a theory of frontal-posterior underconnectivity. *Neuroscience & Biobehavioral Reviews*, 36(4):1292–1313, 2012.
- [9] Marjane Khodatars, Afshin Shoeibi, Delaram Sadeghi, Navid Ghaasemi, Mahboobeh Jafari, Parisa Moridian, Ali Khadem, Roohallah Alizadehsani, Assef Zare, Yinan Kong, Abbas Khosravi, Saeid Nahavandi, Sadiq Hussain, U. Rajendra Acharya, and Michael Berk. Deep learning for neuroimaging-based diagnosis and rehabilitation of autism spectrum disorder: A review. *Computers in Biology and Medicine*, 139:104949, 2021.
- [10] Beatriz Remeseiro and Veronica Bolon-Canedo. A review of feature selection methods in medical applications. *Computers in Biology and Medicine*, 112:103375, 2019.
- [11] Rizgar Zebari, Adnan Abdulazeez, Diyar Zeebaree, Dilovan Zebari, and Jwan Saeed. A comprehensive review of dimensionality reduction techniques for feature selection and feature extraction. *Journal of Applied Science and Technology Trends*, 1(2):56 – 70, May 2020.
- [12] Xinyu Guo, Kelli C. Dominick, Ali A. Minai, Hailong Li, Craig A. Erickson, and Long J. Lu. Diagnosing autism spectrum disorder from brain resting-state functional connectivity patterns using a deep neural network with a novel feature selection method. *Frontiers in Neuroscience*, 11, 2017.
- [13] Canhua Wang, Zhiyong Xiao, and Jianhua Wu. Functional connectivity-based classification of autism and control using svm-rfcv on rs-fmri data. *Physica Medica*, 65:99–105, 2019.
- [14] Zhi-An Huang, Zexuan Zhu, Chuen Heung Yau, and Kay Chen Tan. Identifying autism spectrum disorder from resting-state fmri using deep belief network. *IEEE Transactions on Neural Networks and Learning Systems*, 32(7):2847–2861, 2021.
- [15] Beibin Li, Erin Barney, Caitlin Hudac, Nicholas Nuechterlein, Pamela Ventola, Linda Shapiro, and Frederick Shic. Selection of eye-tracking stimuli for prediction by sparsely grouped input variables for neural networks: Towards biomarker refinement for autism. In *ACM Symposium on Eye Tracking Research and Applications*, pages 1–8, 2020.
- [16] Tomasz Latkowski and Stanislaw Osowski. Data mining for feature selection in gene expression autism data. *Expert Systems with Applications*, 42(2):864–872, 2015.
- [17] Peter Washington, Kelley Marie Paskov, Haik Kalantarian, Nathaniel Stockham, Catalin Voss, Aaron Kline, Ritik Patnaik, Brianna Chrisman, Maya Varma, Qandeel Tariq, et al. Feature selection and dimension reduction of social autism data. In *Pacific Symposium on Biocomputing 2020*, pages 707–718. World Scientific, 2019.
- [18] A. Di Martino, C. G. Yan, Q. Li, E. Denio, F. X. Castellanos, K. Alaerts, J. S. Anderson, M. Assaf, S. Y. Bookheimer, M. Dapretto, B. Deen, S. Delmonte, I. Dinstein, B. Ertl-Wagner, D. A. Fair, L. Gallagher, D. P. Kennedy, C. L. Keown, C. Keyzers, J. E. Lainhart, C. Lord, B. Luna, V. Menon, N. J. Minshew, C. S. Monk, S. Mueller, R. A. Müller, M. B. Nebel, J. T. Nigg, K. O’Hearn, K. A. Pelphrey, S. J. Peltier, J. D. Rudie, S. Sunaert, M. Thioux, J. M. Tyszka, L. Q. Uddin, J. S. Verhoeven, N. Wenderoth, J. L. Wiggins, S. H. Mostofsky, and M. P. Milham. The autism brain imaging data exchange: towards a large-scale evaluation of the intrinsic brain architecture in autism. *Molecular psychiatry*, 19(6):659–667, 2014. 23774715[pmid] PMC4162310[pmcid] mp201378[PII].
- [19] Amirali Kazeminejad and Roberto C. Sotero. Topological properties of resting-state fmri functional networks improve machine learning-based autism classification. *Frontiers in Neuroscience*, 12:1018, 2019.
- [20] Alexandre Abraham, Michael P. Milham, Adriana Di Martino, R. Cameron Craddock, Dimitris Samaras, Bertrand Thirion, and Gael Varoquaux. Deriving reproducible biomarkers from multi-site resting-state data: An autism-based example. *NeuroImage*, 147:736–745, 2017.
- [21] Anibal Sólton Heinsfeld, Alexandre Rosa Franco, R. Cameron Craddock, Augusto Buchweitz, and Felipe Meneguzzi. Identification of autism spectrum disorder using deep learning and the abide dataset. *NeuroImage: Clinical*, 17:16–23, 2018.
- [22] Yazhou Kong, Jianliang Gao, Yunpei Xu, Yi Pan, Jianxin Wang, and Jin Liu. Classification of

- autism spectrum disorder by combining brain connectivity and deep neural network classifier. *Neurocomputing*, 324:63–68, 2019. Deep Learning for Biological/Clinical Data.
- [23] Zeinab Sherkatghanad, Mohammadsadegh Akhondzadeh, Soorena Salari, Mariam Zomorodi-Moghadam, Moloud Abdar, U. Rajendra Acharya, Reza Khosrowabadi, and Vahid Salari. Automated detection of autism spectrum disorder using a convolutional neural network. *Frontiers in Neuroscience*, 13:1325, 2020.
- [24] Siddharth Shrivastava, Upasana Mishra, Nitisha Singh, Anjali Chandra, and Shrish Verma. Control or autism - classification using convolutional neural networks on functional mri. In *2020 11th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, pages 1–6, 2020.
- [25] Zhicheng Jiao, Hongming Li, and Yong Fan. Improving diagnosis of autism spectrum disorder and disentangling its heterogeneous functional connectivity patterns using capsule networks. In *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*, pages 1331–1334, 2020.
- [26] Taban Eslami, Vahid Mirjalili, Alvis Fong, Angela R. Laird, and Fahad Saeed. Asd-diagnet: A hybrid learning approach for detection of autism spectrum disorder using fmri data. *Frontiers in Neuroinformatics*, 13:70, 2019.
- [27] Cameron Craddock, Benhajali Yassine, Chu Carlton, Chouinard Francois, Evans Alan, Jakab András, Khundrakpam Budhachandra, Lewis John, Li Qingyang, Milham Michael, Yan Chaogan, and Bellec Pierre. The neuro bureau preprocessing initiative: open sharing of preprocessed neuroimaging data and derivatives. *Frontiers in neuroinformatics.*, 7, 2013.
- [28] N. Zhang X. Yang, P.T. Schrader. A deep neural network study of the abide repository on autism spectrum classification. *International Journal of Advanced Computer Science and Applications*, 11(4):1–6, 2020.
- [29] R. Cameron Craddock, G.Andrew James, Paul E. Holtzheimer III, Xiaoping P. Hu, and Helen S. Mayberg. A whole brain fmri atlas generated via spatially constrained spectral clustering. *Human Brain Mapping*, 33(8):1914–1928, 2012.
- [30] Xin Yang, Mohammad Samiul Islam, and A M Arefin Khaled. Functional connectivity magnetic resonance imaging classification of autism spectrum disorder using the multisite abide dataset. In *2019 IEEE EMBS International Conference on Biomedical Health Informatics (BHI)*, pages 1–4, 2019.
- [31] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2014.
- [32] Carl Doersch. Tutorial on variational autoencoders, 2021.
- [33] Nitish Srivastava Geoffrey Hinton and Kevin Swersky. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent, 2012.
- [34] Kenji Matsuura Cort J. Willmott. Advantages of the mean absolute error (mae) over the root mean square error (rmse) in assessing average model performance. *Climate Research*, 30:79–82, 2005.
- [35] Yi-Wei Chen and Chih-Jen Lin. *Combining SVMs with Various Feature Selection Strategies*, pages 315–324. Springer Berlin Heidelberg, Berlin, Heidelberg, 2006.