

FaFCNN: A General Disease Classification Framework Based on Feature Fusion Neural Networks

Menglin Kong^a, Shaojie Zhao^b, Juan Cheng^a, Xingquan Li^{c,d}, Ri Su^a, Muzhou Hou^a, Cong Cao^{a*}

^a School of Mathematics and Statistics, Central South University, Changsha, China

^b School of Mathematics, Physics & Statistics, Shanghai University of Engineering Science, Shanghai, China

^c Peng Cheng Laboratory, Shenzhen, China

^d School of Mathematics and Statistics, Minnan Normal University, Zhangzhou, China

{212112025, suricsu, hmzw, congcao}@csu.edu.cn, chengjuan0306@163.com, m440121303@sues.edu.cn, fzulxq@gmail.com

Abstract—There are two fundamental problems in applying deep learning/machine learning methods to disease classification tasks, one is the insufficient number and poor quality of training samples; another one is how to effectively fuse multiple source features and thus train robust classification models. To address these problems, inspired by the process of human learning knowledge, we propose the Feature-aware Fusion Correlation Neural Network (FaFCNN), which introduces a feature-aware interaction module and a feature alignment module based on domain adversarial learning. This is a general framework for disease classification, and FaFCNN improves the way existing methods obtain sample correlation features. The experimental results show that training using augmented features obtained by pre-training gradient boosting decision tree yields more performance gains than random-forest based methods. On the low-quality dataset with a large amount of missing data in our setup, FaFCNN obtains a consistently optimal performance compared to competitive baselines. In addition, extensive experiments demonstrate the robustness of the proposed method and the effectiveness of each component of the model¹.

Index Terms—Disease classification, Neural networks, Feature fusion, Domain adversarial learning

I. INTRODUCTION

With the ability to use large amounts of precisely labelled data, deep neural network-based approaches have achieved exciting results for tasks such as e-commerce recommendation systems [1], image classification [2], and object detection [3]. However, many tasks in the medical field tend to have insufficient samples and a large amount of missing data, which makes it extremely difficult to develop a general deep-learning framework for disease classification tasks in the medical field [4] [5]. In addition, the records corresponding to patients in hospital databases often involve demographic features, clinical features, radiological features, and other diagnostic metrics from multiple sources; there are often scale inconsistencies and information redundancy among these data, and using them together to train machine learning models may compromise the interpretability and robustness

of the models [6]. In summary, there are two fundamental problems in applying deep learning/machine learning methods to disease classification tasks: (1) the insufficient number and poor quality of training samples; (2) how to effectively fuse multiple source features and thus train robust classification models.

To address these problems, inspired by the process of human learning knowledge, some researchers [7]–[9] propose to augment the feature representation of each sample using the features of similar samples in the training set, i.e., introducing sample correlation features as an extension of existing features. As in [7], the authors make an adjustment to the coupled two-stage modelling by directly using the prediction probabilities of the random forest (RF) model as correlation features with the original features as input, using a DNN with a two-tower structure to map the two parts of features separately, and finally making predictions based on the summation of high-level features. However, the prediction probability of the RF model of each sample is not enough to characterize the similarity with other samples in the training set, which will impair the performance of the model. The literature [8] proposes a graph generation method for medical datasets based on sample paths of a pre-trained random forest (RF) model, transforming structured data into graph data and training a graph convolutional network for node classification to achieve accurate differentiation of Crohn’s disease and intestinal tuberculosis. Nevertheless, this method relies heavily on artificial thresholds to determine the edges between nodes when constructing graph data, which leads to poor robustness of the framework.

AI has shown powerful potential in the field of data-driven medical fields. Esteva et al. [4] elaborated on the application prospects of various methods in the field of deep learning in the medical field from four aspects: computer vision, natural language processing, reinforcement learning and generalized deep learning methods. Rauschert et al. [10] briefly summarized the current state of Machine Learning (ML), and showed that recent advances in deep learning offer greater promise in helping physicians achieve accurate diag-

Cong Cao is the corresponding author (congcao@csu.edu.cn).

¹Accepted in IEEE SMC2023

noses. For example, Lima et al. proposed FSTBSVM [11], a twin-bounded SVM classifier combined with a scalable feature selection method. And then, Kuma et al. proposed a classification algorithm that combines k -nearest neighbour and genetic algorithm [12]; Gu et al. proposed a fuzzy support machine with the Gaussian kernel as well as linear kernel [13]. However, due to the problems of small sample size and incomplete data in medical datasets, existing studies basically design special classification algorithms for specific disease classification tasks and basically follow the paradigm of feature selection plus machine learning model prediction. At present, there is no unified and generalized framework for auxiliary diagnosis of medical diseases.

Considering the advantages and disadvantages of existing methods, we propose the **Feature-aware Fusion Correlation Neural Network (FaFCNN)**, a general framework for disease classification. Specifically, we keep the idea of using an agent model to obtain sample correlation features to realize feature augmentation from existing methods, while the sample correlation features are acquired based on the positions of samples in the leaf nodes of a pre-trained gradient boosting decision tree (GBDT). It is experimentally demonstrated that the augmented features obtained by our method capture more accurate sample correlation than the RF-based augmented features, and further improve the model performance. In order to further improve the performance of disease classification models on low-quality datasets, FaFCNN considers the correlation of features in addition to the correlation of samples, and introduces a feature-aware interaction module (FaIM) and a feature alignment module (FAM) based on domain adversarial learning to achieve more efficient feature fusion and model performance.

The contributions of this paper are listed below:

- We propose FaFCNN, a generic deep learning-based framework for disease classification, and our method obtains a consistently optimal performance compared to competitive baselines.
- We improve the way existing methods obtain sample correlation features, training using augmented features obtained by pre-training GBDT yields more performance gains than RF-based methods.
- In the feature fusion approach, the feature alignment module based on domain adversarial learning introduced by FaFCNN alleviates the performance degradation caused by the naive summation of existing methods.
- We synthesise low-quality datasets by adding different levels of perturbation on four public datasets. Extensive experiments demonstrate the robustness of our proposed method and the effectiveness of each component of the model.

II. METHODOLOGY

A. Correlation Features Construction

In this section, we present the construction of sample similarity features based on pre-trained GBDT. GBDT [14] [15] is an integrated model consisting of decision trees that learn in a gradient-boosting manner, where each base

classifier (DT) is trained to fit the residuals of the prediction results of the preorder model. The GBDT structure is shown in “Fig. 1”.

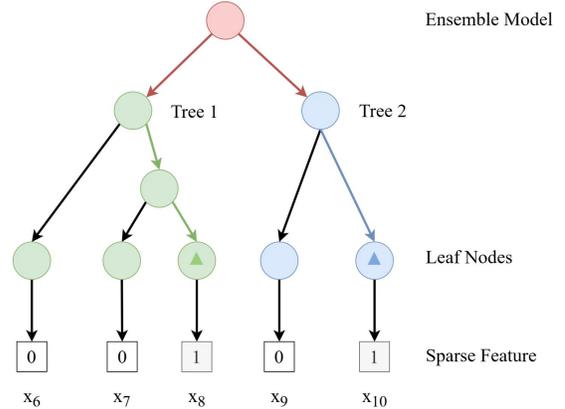


Fig. 1. The diagram of correlation features construction based on GBDT. The red circle represents the root node, the green circles represents the middle node and leaf node of the first base classifier, the blue circle represents the second base classifier, and the triangles represent the position of the sample to be predicted in the leaf node of the base classifier.

In a more general scenario, the sample correlation feature construction method can be expressed as follows: first train a GBDT model on the full training data D_{train} with the number of base classifiers M and the number of leaf nodes of each base classifier k . For a sample $\mathbf{x} \in \mathbb{R}^d$ with d -dimensional features, it is fed into the model to get the prediction result and its position in the leaf nodes is recorded and a $k \times M$ -dimensional one-hot vector is obtained, i.e., $\mathbf{x}_{aug} = (0, 1, 0, \dots, 0) \in \mathbb{R}^{k \times M}$, and finally the \mathbf{x} is concatenated with the original features \mathbf{x} to the augmented feature vector $\tilde{\mathbf{x}} \in \mathbb{R}^{k \times M + d}$. In this way, the generated one-hot vector represents the position of the sample’s leaf nodes in GBDT, i.e., the predicted path of the sample given by each base classifier. From the perspective of partitioning the feature space, the prediction path of each base classifier corresponds to the subregion in the feature space where the sample points are located in a certain view. The more the intersection of the prediction paths of two samples, the more they are in the same subregion in multiple views, i.e., they have a higher correlation. We can use this correlation to augment the sample features to train a deep neural network with more powerful representation ability.

B. Feature-aware Interaction Module

FaFCNN improves the naive FCNN in two ways, respectively, by introducing FaIM to perform correlation-based mapping (i.e., feature interaction) on the sample correlation features \mathbf{x}_{aug} , which results in a finer-grained intermediate representation about \mathbf{x}_{aug} ; and by introducing FAM to perform domain adversarial learning-based feature alignment operations on the original features \mathbf{x} and the intermediate representation of the sample correlation features \mathbf{x}_{aug} .

In this section, we detail how FaIM obtains fine-grained intermediate representations about sample correlation features \mathbf{x}_{aug} by modeling feature interaction. As illustrated in

the purple box of “Fig. 2”, considering a sample with 5-dimensional sample correlation feature $\mathbf{x}_{\text{aug}} = (1, 0, 1, 0, 1)$, we initialize a p -dimensional vector \mathbf{h}_i for the i th dimension in \mathbf{x}_{aug} to obtain 5 p -dimensional vectors, each p -dimensional vector characterizes richer semantic information of each dimension of \mathbf{x}_{aug} . Then the vector h_i , where $i \in \{i|x_i = 1\}$ corresponding to the non-zero position is taken to perform the second-order interactions between features in an element-wise product manner. Attention Net, a sub-network with softmax activation function, calculates the weight $a_{i,j}$ for each feature interaction term $h_i \odot h_j$ (where $(i, j) \in \{m|x_m = 1\}$) in a self-attention manner, and finally uses this weight to aggregate these second-order interaction features to obtain the mapped sample correlation features. More generally, consider a sample of $\mathbf{x}_{\text{aug}} \in \mathbb{R}^{k \times M}$, the p -dimensional vector \mathbf{h}_{aug} after being mapped can be obtained by the following formula:

$$\mathbf{h}_{\text{aug}} = \sum_{i=1}^{k \times M} w_i x_i + \sum_{i=1}^{k \times M} \sum_{j=i+1}^{k \times M} a_{ij} (\mathbf{h}_i \odot \mathbf{h}_j) x_i x_j \quad (1)$$

where the weight $a_{i,j}$ is calculated by the following formula:

$$\begin{aligned} a'_{ij} &= \mathbf{q}^T \text{ReLU}(\varpi_{\text{attn}} (\mathbf{h}_i \odot \mathbf{h}_j) x_i x_j + b_{\text{attn}}), \\ a_{ij} &= \frac{\exp(a'_{ij})}{\sum_{(i,j) \in \mathcal{I}_{\mathbf{x}_{\text{aug}}}} \exp(a'_{ij})} \end{aligned} \quad (2)$$

where $\mathbf{q}, \varpi_{\text{attn}}, b_{\text{attn}}$ are the parameters of the sub-network, $\mathcal{I}_{\mathbf{x}_{\text{aug}}}$ denotes the set of the index in \mathbf{x}_{aug} .

Due to the huge number of two-by-two combinations between features, for example, for $k \times M$ -dimensional features one needs to compute $C_2^{k \times M} = (k \times M) \times (k \times M - 1) / 2$ feature interaction terms and their weights, from the perspective of enhancing the interpretability of the model and reducing the computation, we want most of the interaction terms to have a weight equals to zero. This not only highlights the combination of features with the greatest impact on the prediction but also greatly reduces the computation. Inspired by the addition of L1-norm-based regularization terms to the coefficients of the linear model in LASSO regression [16], FaFCNN adds L1-norm-based sparse regularization terms to the output $a_{i,j}$ of Attention Net in the expectation of compressing the weights of unimportant feature combinations toward the value of zero and highlighting the important ones, the formula is as follows:

$$L_{\text{sparse}} = \sum_{i=1}^{k \times M} \sum_{j=i+1}^{k \times M} \|a_{ij}\|_1 \quad (3)$$

C. Feature Alignment Module

FaFCNN introduces adversarial-learning-based FAM to achieve a smoother feature fusion by aligning the distribution of mapped original features \mathbf{x} and the sample correlation features \mathbf{x}_{aug} in the high-dimensional representation space, which is shown in the orange box of Fig.2(a).

Similar to FCNN, a neural network with two hidden layers (DNN in Fig.2(a)) is first used to map the original features

to obtain their representations in high-dimensional space $\mathbf{h} \in \mathbb{R}^p$, the formula is as follows:

$$\mathbf{h} = f_{o,2}(\varpi_{o,2} \cdot f_{o,1}(\varpi_{o,1} \cdot \mathbf{x} + b_{o,1}) + b_{o,2}) \quad (4)$$

Where $\varpi_{o,1}, \varpi_{o,2}, b_{o,1}, b_{o,2}$ are the parameters of DNN. Since FaFCNN uses different mapping methods for different features (correlation-based aggregation for \mathbf{x}_{aug} , and MLP-based nonlinear mapping for \mathbf{x}), which leads to a large difference in the distribution of the two parts of features in the high-dimensional representation. Therefore, FaFCNN introduces the FAM module to achieve the distribution alignment of \mathbf{h} and \mathbf{h}_{aug} in the representation space with the idea of Min-Max game in generative adversarial networks (GAN) [17].

FaFCNN introduces a discriminator D in FAM to distinguish signals from two partial features with \mathbf{h} and \mathbf{h}_{aug} as inputs, and the optimization objective is to enhance the discriminator’s ability to distinguish \mathbf{h} and \mathbf{h}_{aug} , i.e:

$$\theta^* = \underset{\theta}{\text{argmax}} \|D(\mathbf{h}; \theta) - D(\mathbf{h}_{\text{aug}}; \theta)\|_1 \quad (5)$$

which is equal to minimizing the following formula:

$$L_D = - \sum_{i=1}^N \|D(\mathbf{h}_i; \theta) - D(\mathbf{h}_{\text{aug},i}; \theta)\|_1 \quad (6)$$

where θ is the parameter of the D , which is a two-layer MLP in FaFCNN. Naturally, we can consider the above-mentioned DNN that maps \mathbf{x} as the generator G in GAN, whose optimization goal is to make the distribution of the mapped \mathbf{h} in the high-dimensional representation space as similar as possible to \mathbf{h}_{aug} , so that the discriminator D cannot distinguish \mathbf{h}_{aug} from \mathbf{h} , the formula is as follows:

$$\phi^* = \underset{\phi}{\text{argmin}} \|D(G(\mathbf{x}; \phi); \theta) - \mathbf{1}\|_1 \quad (7)$$

where $\phi = \{\varpi_{o,1}, \varpi_{o,2}, b_{o,1}, b_{o,2}\}$ is the parameter of the DNN, θ^* is the optimal parameters of the discriminator in the last iteration, $\mathbf{1} \in \mathbb{R}^p$ is an all-one vector(proxy label of \mathbf{h}_{aug} in this domain adversarial learning procedure). This optimal ϕ^* can be found by minimizing the following loss function:

$$L_G = \sum_{i=1}^N \|D(G(\mathbf{x}_i; \phi); \theta) - \mathbf{1}\|_1 \quad (8)$$

However, it is easy to experience pattern collapse during training with GAN, i.e., the generator generates very narrow distributions that cover only a single pattern in the data distribution. This was also observed in our experiments, where DNNs tend to consistently map samples with different original features \mathbf{x} to limited range in the high-dimensional space during adversarial learning, but since this is a pattern in the \mathbf{h}_{aug} distribution, the purpose of tricking the discriminator D can be reached. While this single-pattern representation does not bring any beneficial information to the model for classification. To ensure that the aligned \mathbf{h} maintains diversity during the adversarial learning, FaFCNN introduces a supervised signal to \mathbf{h} so that it retains a certain

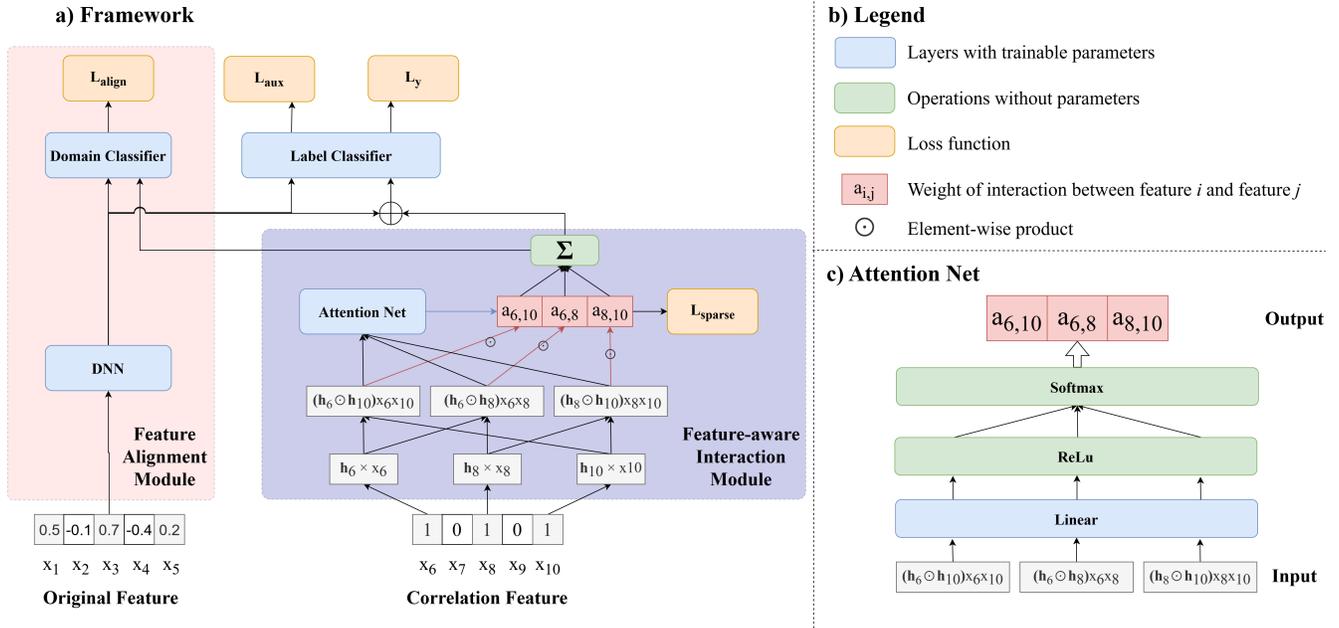


Fig. 2. The structural diagram of the proposed FaFCNN. (a) The overall framework of FaFCNN. The part in the orange box is FAM, the part in the purple box is FaIM. (b) is an explanation of those graphics that appear previously. (c) The forward calculation process in the Attention Net of the FaIM

amount of information that is beneficial to the classification of the sample thus ensuring that the aligned \mathbf{h} diversity of the distribution pattern, noting the label classifier as $F(\cdot; \psi)$ and the auxiliary loss as follows:

$$L_{aux} = -\frac{1}{N} \sum_{i=1}^N (y_i \log F(\mathbf{h}_i; \psi) + (1 - y_i) \log (1 - F(\mathbf{h}_i; \psi))) \quad (9)$$

D. Optimization

Based on the above, the training process of FaFCNN consists of two stages. In the first stage, we train the \mathbf{h}_{aug} obtained from FaIM to make it capable of classifying samples in a supervised manner, while adding a sparse regularization term to the total loss to ensure the sparsity of the weights of feature interaction terms obtained from Attention Net, the formula is as follows:

$$L_y = -\frac{1}{N} \sum_{i=1}^N (y_i \log F(\mathbf{h}_i, \mathbf{aug}; \psi) + (1 - y_i) \log (1 - F(\mathbf{h}_i, \mathbf{aug}; \psi))) \quad (10)$$

$$L_1 = L_y + \alpha L_{sparse} \quad (11)$$

In the second stage, we first freeze the network parameters in the already trained FaIM module to ensure that \mathbf{h}_{aug} does not change during FAM training. Then alternately optimize the parameters θ of the discriminator D with Equation 6, and the parameters ϕ of the DNN with Equation (8) and (9), the formula is as follows:

$$L_2 = L_{aux} + \beta L_G \quad (12)$$

III. EXPERIMENTS

In this section, We validate the effectiveness and robustness of FaFCNN on four publicly available medical datasets with special perturbation treatments.

A. Experimental Setting.

1) *Dataset*: To prove the superiority of the proposed method in medical diagnosis, we apply our model to four public medical datasets, including the Wisconsin Breast cancer, Pima Indians Diabetes, Hepatitis, Heart-Statlog datasets, more details of these datasets are listed as follows:

TABLE I
PUBLIC DATASETS DESCRIPTION

Dataset	N-smamples	N-features	N-classes	Reference
Wisconsin Breast Cancer	699	9	2	[18]
Pima Indians Diabetes	768	8	2	[19]
Hepatitis	155	19	2	[20]
Heart-Statlog	270	13	2	[21]

To simulate the challenge of existing a large number of missing values in a real scenario medical dataset, we add different levels of perturbation to the above dataset. Considering a raw dataset with N samples and d features, the data preprocessing process is as follows:

- First, the missing values in the dataset are processed. The columns with missing values are first identified and the median of the column other than the missing values is calculated and the missing values are replaced by the median.
- Then, the dataset is perturbed randomly. The data are first shuffled by rows, and then the rows of data to be

perturbed are removed according to the selected δ . Each column of these data rows is randomly selected with equal probability ($1/d$) and perturbed in the same way as the missing values are processed.

- Finally, the data set is divided into a training set, validation set, and test set in the ratio of 8:1:1.

2) *Hyperparameter*: In this section the hyperparameter settings used for training FaFCNN are described, the parameters of the GBDT included $k = \text{integer}(d/2)$ (d is the number of features of the dataset) estimators, 8 max depths, and $M = 2$ min sample leaves. The stage 1 training of FaFCNN uses the Adam optimizer with a learning rate of 0.005 for $T_1 = 10000$ epochs. The stage 2 training of DNN and Discriminator uses the SGD optimizer with a learning rate of 0.005 for $T_2 = 10000$ epochs. The dimension p of vector \mathbf{h} is set as 8, and the balance coefficient α and β are set to 0.05 and 0.5, respectively.

3) *Performance evaluation metrics*: We select accuracy, sensitivity and Specificity, which are common evaluation metrics in classification tasks, to compare the performance of FaFCNN and baseline models in three dimensions. The three evaluation indicators are defined as follows:

$$Acc = \frac{\mathcal{M}_{tp} + \mathcal{M}_{tn}}{\mathcal{M}_{tp} + \mathcal{M}_{fp} + \mathcal{M}_{tn} + \mathcal{M}_{fn}} \quad (13a)$$

$$Sensitivity = \frac{\mathcal{M}_{tp}}{\mathcal{M}_{tp} + \mathcal{M}_{fn}} \quad (13b)$$

$$Specificity = \frac{\mathcal{M}_{tn}}{\mathcal{M}_{tn} + \mathcal{M}_{fp}} \quad (13c)$$

B. Classification results of different classification methods

Our comparison is performed uniformly on four datasets with a perturbation ratio of $\delta = 0.5$; to ensure the fairness of the comparison, the structure of the DL-based baseline is adjusted so that the number of parameters of the models involved in the comparison remains the same. The experimental results are shown in Table II, and results in the table are the mean values of 10 independent repetitions of the experiment.

As shown in Table II, in the comparison of results from the Wisconsin Breast Cancer dataset, the DL-based methods (DNN, RFG-GCN) do not show better performance in some metrics than the ML-based methods (RF, LR); the FCNN better exploits the sample correlation in the training set to achieve a consistent performance improvement over the DNN. FaFCNN achieves smoother and more effective feature fusion while considering feature correlation and achieves significant improvement in two evaluation metrics compared to FCNN.

ML-based methods perform poorly on the Pima Indians Diabetes dataset, especially the sensitivity metric does not exceed 70% at the highest; DL-based methods achieve improvements in the other two metrics, but the sensitivity metric still can not exceed 80% (78.1% for FCNN). This indicates that there is a serious class imbalance problem on the Diabetes dataset and the model easily misclassifies some of the positive cases as negative cases. FaFCNN significantly

outperforms FCNN with a p-value of 0.001 at 91.5% on the sensitivity metric while the other two metrics significantly outperform the optimal baseline with a p-value of 0.005.

On the Hepatitis dataset, the SRLPSO-ELM method achieves optimal performance in accuracy but this advantage was not significant (98.7%) and our FaFCNN has achieved (98.6%); FSTBSVM reaches 100% in specificity but does not significantly outperform FaFCNN (98.5%), and this advantage comes at the expense of sensitivity (78.6%), while FaFCNN achieves optimal performance in this metric (98.7%).

On the Heart-Statlog dataset, DL-based methods consistently outperform ML-based methods, and FaFCNN again achieves optimal performance on the three metrics with the same number of parameters and demonstrates significance in terms of accuracy and sensitivity due to the well-designed structure of the network. In summary, our FaFCNN is able to show robust and consistent optimal performance with respect to the baseline models on multiple datasets with 50% of the samples perturbed in a low-quality data setting and with class imbalance problems.

C. Robustness verification of classification results

To verify that our proposed FaFCNN maintains robustness and acceptable performance in the face of scenarios with large amounts of missing data, we set a set of perturbation ratios $\delta \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$, 10 experiments are conducted on the Wisconsin Breast Cancer dataset for each δ , and the results of the experiments are shown in “Fig. 3”.

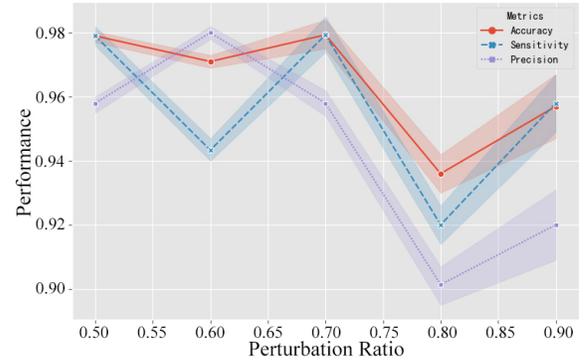


Fig. 3. Performance of FaFCNN on the Wisconsin Breast Cancer dataset with different settings of perturbation ratio δ . The red solid line, blue dashed line, and purple dotted line represent accuracy, sensitivity, and precision, respectively, and the points on the axes represent the mean values of 10 experiments, while the upper and lower bandwidths represent the standard deviation of the experimental results.

As shown in “Fig. 3”, we can conclude that as δ gradually increases, meaning that the proportion of samples with missing values in the dataset increases, the fluctuation of FaFCNN’s performance also gradually increases (the bandwidth on both sides of the performance line increases), but the three evaluation metrics still maintain a high level (the mean value of the worst case also remains above 0.9). Specifically, accuracy does not decrease significantly as δ increases, and the mean value of each case remains above 0.93. Sensitivity and Precision show a decreasing trend (when δ increases from

TABLE II
PERFORMANCE RESULTS OF DIFFERENT MODELS ON LARGE COMPETITION DATA SETS

Methods		Datasets											
		Wisconsin Breast Cancer			Pima Indians Diabetes			Hepatitis			Heart-Statlog		
		Accuracy(%)	Sensitivity(%)	Specificity(%)	Accuracy(%)	Sensitivity(%)	Specificity(%)	Accuracy(%)	Sensitivity(%)	Specificity(%)	Accuracy(%)	Sensitivity(%)	Specificity(%)
ML-based	LR	90.7	79.5	86.1	76.6	51.8	76.3	83.5	67.2	92.5	87	71.4	93.8
	RF	92.1	91.5	86	76.6	49	68.6	86.7	85.5	71	83.3	71.4	83.3
	PSO-ELM [8]	93.6	89.8	91.7	78.6	56.4	77.5	97.4	93.7	95.7	85.9	86	86
	SRLPSO-ELM [8]	91.4	84.6	84.6	74	50	65	98.7	94.2	96	89.9	87.8	88.4
	ESTBSVM [11]	95	92.3	90	76.6	69.2	64.3	89.1	78.6	100	85.9	72.7	89.8
	KNN-GA [12]	90	84.8	84.8	72.1	33.9	83.3	64.4	78.3	50	83.5	87.5	82.3
	FSVM-Gaussian [13]	92.9	90.7	86.7	76	52.2	61.5	84	89.6	87.4	84.7	86	83.1
DL-based	RF+GCN [8]	90	80.9	88.4	83.5	67.2	93.5	95.9	<u>98.1</u>	92.2	92.1	90.2	93.5
	DNN [7]	90	85.1	85.1	84.4	69.2	93.4	91	93.8	89.1	92.6	89.7	<u>94.8</u>
	FCNN [7]	<u>95.7</u>	<u>97.5</u>	<u>88.6</u>	<u>88.3</u>	<u>78.1</u>	<u>93.9</u>	90.8	95.7	87.3	<u>95.1</u>	<u>91.3</u>	93.5
	FaFCNN(ours)	97.9*	97.9	95.8*	93.1*	91.5**	96.3*	98.6	98.7	98.5	97.7*	94.9*	95.2

Accuracy comparison of ours and baseline models on four 50% perturbed datasets. The best is in bold font, * and ** indicate $p < 0.005$ and $p < 0.001$ for a one-tailed t-test, and the underlining indicates sub-optimal performance. The results are averaged over 10 independent repeats on all the datasets.

0.7 to 0.8) and show a performance increase when δ improves to 0.9. In summary, FaFCNN has a narrow bandwidth on both sides of the line for different values of δ , which proves the robustness of the classification results in each setting; its classification performance does not show an obviously decreasing trend with the increase of δ , which proves that the model has strong robustness to noisy data.

D. Ablation Study

In this section, we focus on validating the effectiveness of the well-designed components in FaFCNN by means of ablation experiments; to ensure fairness of the comparison, each variant of FaFCNN is extended in terms of network structure to ensure consistent overall model parameters. We conduct 10 independent repetition experiments on the Wisconsin Breast Cancer dataset under the setting of $\delta = 0.5$.

1) *Validity of FaIM&FAM*: We first validate the effectiveness of the proposed modules and quantify the performance improvement brought by each module through a set of comparison experiments between FaFCNN and its three variants on the Wisconsin Breast Cancer dataset, as shown in “Fig. 4”.

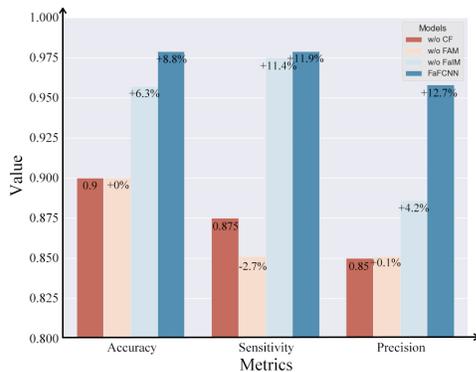


Fig. 4. Comparative results of FaFCNN and its three variants on Wisconsin Breast Cancer. Dark red means no sample correlation features are used, light orange means sample correlation features are added but FAM is not used, light blue means correlation between features is not modelled using FaIM, and dark blue means FaFCNN. The number above the bar indicates the relative improvement of adding different modules compared to the base model without introducing sample correlation.

As shown in “Fig. 4”, compared with the base model without sample correlation features, using the output of pre-trained RF as augmented features do not improve the

performance of the model, but decreases the sensitivity of the model (-2.74%), which indicates that the introduction of augmented features without using a reasonable feature fusion method will harm the performance of the model. The introduction of FAM significantly improves the performance of the model, with improvements of 6.3%, 11.4%, and 4.2% for the three metrics respectively, which validates the effectiveness of our proposed feature fusion module based on adversarial learning. FaFCNN uses the FaIM module to replace the w/o FaIM variant’s DNN for mapping sample correlation features and achieves further improvements in accuracy(+8.8%), sensitivity(+11.9%) and precision(+12.7%) metrics, which verifies that using predicted paths of samples of GBDT as augmented features can capture more accurate sample correlation than the RF-based approach, and the feature-interaction-based explicit mapping approach can achieve finer-grained feature representation than the DNN-based implicit feature mapping.

2) *Effectiveness of Sparse Regularization*: To verify the effectiveness of the weight sparse regularization term added to the FaIM module, we train FaFCNN and FaFCNN without sparse regularization on the 50% perturbed Wisconsin Breast Cancer dataset, respectively, and record the output of the Attention Network, i.e., the weights of the feature interactions for each sample in the test phase, then average on them. “Fig. 5” shows the heatmap based on the mean value of 10 repetitions of the above procedure.

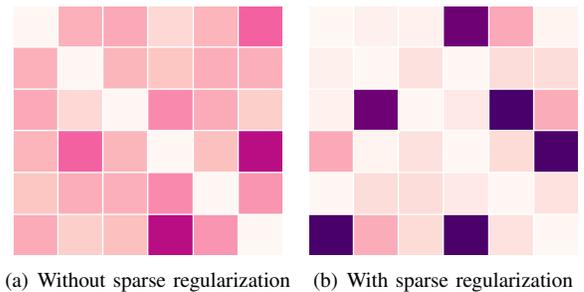


Fig. 5. The heatmap of average weights of feature interactions in FaIM, calculated in the test phase of 50% perturbed Wisconsin Breast Cancer dataset. The darker the color, the greater the absolute value of the weight.

FaFCNN and FaFCNN-w/o sparse regularization perform consistently on 10 independent replicate experiments, that us-

ing sparse regularization with mean values of 97.9%, 97.9%, 95.8% and Without sparse regularization with mean values of 95%, 92%, 93.9% for the three metrics, respectively, and FaFCNN shows a significant improvement in accuracy and sensitivity relative to the variant without sparse regularization term. On the other hand, the above heat map shows that the two models capture similar feature association patterns, such as the larger values of $a_{3,5}$, which implies that the feature interaction terms at these two locations have a greater impact on the model prediction thus the two features are more correlated. In addition, the sparse regularization term does work as expected by reducing the weights of relatively unimportant feature interactions while increasing the weights of critical feature interactions (Fig.5(b) has many more blank squares than Fig.5(a) but with darker colors at important positions), allowing the model to discover significant feature interaction patterns in the data, thus reducing the computation by using only the important feature combinations in the subsequent modelling process.

IV. CONCLUSIONS

In this work, by considering the advantages and disadvantages of existing methods, we propose the FaFCNN, a general framework for disease classification. On the one hand, FaFCNN improves the way existing methods obtain sample correlation features, exploiting augmented features obtained by pre-training gradient boosting decision trees to capture more accurate correlations between samples in the training set. On the other hand, FaFCNN introduces a feature alignment module for smoother and more efficient feature fusion, and the feature-aware interaction module considers feature correlation and model feature interaction in a more fine-grained manner to enhance the model's representation ability. Extensive experimental results show that FaFCNN has strong robustness and can achieve consistent optimal performance concerning the baseline models on multiple datasets with 50% of the samples perturbed in a low-quality data setting and with class imbalance problems.

Acknowledgements This study was supported by Natural Science Foundation of Hunan Province of China (grant number 2022JJ30673) and by the Graduate Innovation Project of Central South University (2023XQLH032, 2023ZZTS0304).

REFERENCES

- [1] X. He and T.S. Chua, "Neural factorization machines for sparse predictive analytics," *Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval*, pp. 355-364, 2017.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770-778, 2016.
- [3] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, pp. 91-99, 2015.
- [4] A. Esteva, A. Robicquet, B. Ramsundar, et al. "A guide to deep learning in healthcare," *Nature Medicine*, vol. 25, no. 1, pp. 24-29, 2019.
- [5] A. Y. Hannun, P. Rajpurkar, M. Haghpanahi, et al. "Cardiologist-level arrhythmia detection and classification in ambulatory electrocardiograms using a deep neural network," *Nature Medicine*, vol. 25, no. 1, pp. 65-69, 2019.

- [6] E. J. Topol, "High-performance medicine: the convergence of human and artificial intelligence," *Nature Medicine*, vol. 25, no. 1, pp. 44-56, 2019.
- [7] Y. Chen, Y. Li, M. Wu, F. Lu, M. Hou, and Y. Yin, "Differentiating Crohn's disease from intestinal tuberculosis using a fusion correlation neural network," *Knowledge-Based Systems*, vol. 244, pp. 108570, 2022.
- [8] K. Cao, Y. Xiao, and M. Hou "Correlation-driven framework based on graph convolutional network for clinical disease classification," *Journal of Statistical Computation and Simulation*, vol. 91, no. 15, pp. 3108-3124, 2022.
- [9] Z. Zhang, B. Chen, S. Xu, G. Chen, and J. Xie, "A novel voting convergent difference neural network for diagnosing breast cancer," *Neurocomputing*, vol. 437, pp. 339-350, 2021.
- [10] S. Rauschert, K. Raubenheimer, P. E. Melton, and R. C. Huang, "Machine learning and clinical epigenetics: A review of challenges for diagnosis and classification," *Clinical Epigenetics*, vol. 12, no. 1, pp. 1-11, 2020.
- [11] M. D. de Lima, J. D. O. R. E. Lima, and R. M. Barbosa, "Medical data set classification using a new feature selection algorithm combined with twin-bounded support vector machine," *Medical Biological Engineering Computing*, vol. 58, no. 3, pp. 519-528, 2020.
- [12] S. G. Kumar S, "Medical dataset classification using k-NN and genetic algorithm," *Advances in Intelligent Systems and Computing*, vol. 556, no. 11, pp. 813-823, 2017.
- [13] X. Gu, T. Ni, and H. Wang, "New fuzzy support vector machine for the class imbalance problem in medical datasets classification," *The Scientific World Journal*, vol. 2014, pp. 1-12, 2014.
- [14] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Annals of Statistics*, vol. 29, no. 5, pp. 1189-1232, 2001.
- [15] J. Friedman, T. Hastie, and R. Tibshirani, "Additive logistic regression: A statistical view of boosting," *The Annals of Statistics*, vol. 28, no. 2, pp. 337-407, 2000.
- [16] T. Hesterberg, N. H. Choi, L. Meier, and C. Fraley, "Least angle and l_1 regression: A review," *Statistics Surveys*, vol. 2, pp. 61-93, 2008.
- [17] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139-144, 2020.
- [18] W. Wolberg, "Breast Cancer Wisconsin (Original)," *UCI Machine Learning Repository*, 1992. <https://doi.org/10.24432/C5HP4Z>.
- [19] M. Kahn, "Diabetes," *UCI Machine Learning Repository*. <https://doi.org/10.24432/C5T59G>.
- [20] "Hepatitis," *UCI Machine Learning Repository*, 1988. <https://doi.org/10.24432/C5Q59J>.
- [21] "Statlog (Heart)," *UCI Machine Learning Repository*. <https://doi.org/10.24432/C57303>.