



UNIVERSITY OF LEEDS

This is a repository copy of *Parkinson's Disease Classification and Clinical Score Regression via United Embedding and Sparse Learning From Longitudinal Data*.

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/172817/>

Version: Accepted Version

Article:

Huang, Z, Lei, H, Chen, G et al. (5 more authors) (2021) Parkinson's Disease Classification and Clinical Score Regression via United Embedding and Sparse Learning From Longitudinal Data. IEEE Transactions on Neural Networks and Learning Systems. ISSN 2162-237X

<https://doi.org/10.1109/tnnls.2021.3052652>

© 2021 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. Uploaded in accordance with the publisher's self-archiving policy.

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

Parkinson's Disease Classification and Clinical Score Regression via United Embedding and Sparse Learning from Longitudinal Data

Zhongwei Huang, Haijun Lei, Guoliang Chen, Alejandro F Frangi, *Fellow, IEEE*, Yanwu Xu, *Senior Member, IEEE*, Ahmed Elazab, Jing Qin, Baiying Lei*, *Senior Member, IEEE*

Abstract—Parkinson's disease (PD) is known as an irreversible neurodegenerative disease that mainly affects the patient's motor system. Early classification and regression of PD are essential to slow down this degenerative process from its onset. In this paper, a novel adaptive unsupervised feature selection approach is proposed by exploiting manifold learning from longitudinal multimodal data. Classification and clinical score prediction are performed jointly to facilitate early PD diagnosis. Specifically, the proposed approach performs united embedding and sparse regression, which can determine the similarity matrices and discriminative features adaptively. Meanwhile, we constrain the similarity matrix among subjects and exploit $\ell_{2,p}$ norm to conduct sparse adaptive control for obtaining the intrinsic information of the multimodal data structure. An effective iterative optimization algorithm is proposed to solve this problem. We perform abundant experiments on the Parkinson's Progression Markers Initiative (PPMI) dataset to verify the validity of the

proposed approach. The results show our approach boosts performance on the classification and clinical score regression of longitudinal data and surpasses the state-of-the-art approaches.

Index Terms—Classification, clinical score prediction, embedding learning, longitudinal multimodal data, Parkinson's disease, sparse regression.

I. INTRODUCTION

Parkinson's disease (PD) is a common neurodegenerative disorder that occurs in the elderly. With the worsening of the condition, it can trigger incidents leading to death. Patients suffering from PD usually not die from the disease but accidents or complications related to the disease. For example, in advanced cases of the disease, difficulty in swallowing can cause PD patients to inhale food into their lungs, resulting in pneumonia or other pulmonary conditions. Patients with PD have four main motor symptoms: muscle rigidity, static tremor, unstable posture, and bradykinesia [1]. Except for these visual symptoms, there are also some concomitant symptoms (e.g., depression, lethargy, olfaction disorder, and cognition impairment [2]). These symptoms are mainly caused by the degeneration of dopaminergic neurons in a region of the brain called the substantia nigra [3]. However, in early PD therapeutic trials, dopaminergic imaging has found that approximately 15 percent of scans are at the normal level, namely scans without evidence of dopaminergic degeneration (SWEDD) [4]. This condition has undoubtedly augmented the difficulties of PD classification. In the early stage of this disease, it may be challenging to know if the symptoms indicate or imitate PD. Because of this, while early PD classification is essential to slow down this degenerative process from its onset, it is a quite challenging task.

Since multimodal data can offer complementary information for the classification of neurodegenerative diseases, these data have played an increasingly significant role and captured widespread attention [5, 6]. For instance, in [7], multimodal data is utilized to raise the classification performance of neurodegenerative disease based on a semi-supervised feature-subject selection approach. Gray matter (GM) in magnetic resonance imaging (MRI) is widely used to obtain information about changes in nerve cells. First eigenvalue (L1) and first eigenvector (V1) of diffusive tensor imaging (DTI) indicate the largest diffusion coefficient and its direction vector, respectively. Therefore, L1 and V1 may be more sensitive to neurodegeneration in the brain. Inspired by the above, we propose to

This work was supported partly by National Natural Science Foundation of China (Nos. 61871274, 61801305, 81571758, and 61950410615), National Natural Science Foundation of Guangdong Province (Nos. 2020A1515010649 and 2019A151511205), (Key) Project of Department of Education of Guangdong Province (No. 2019KZDZX1015), Guangdong Laboratory of Artificial-Intelligence and Cyber-Economics (SZ), Shenzhen Peacock Plan (Nos. KQTD2016053112051497 and KQTD2015033016104926), Shenzhen Key Basic Research Project (Nos. JCYJ20190808165209410, 20190808145011259, JCYJ20180507184647636, GJHZ20190822095414576, JCYJ20170302153337765, JCYJ20170302150411789, JCYJ20170302142515949, GCZX2017040715180580, GJHZ20180418190529516, and JSGG20180507183215520), NTUT-SZU Joint Research Program (No.2020003), Hong Kong Research Grants Council (No. PolyU 152035/17E), AFF is supported by the RAEng Chair in Emerging Technologies (CIET1819/19).

Z. Huang, H. Lei, and G. Chen are with the Key Laboratory of Service Computing and Applications, Guangdong Province Key Laboratory of Popular High Performance Computers, Guangdong Province Engineering Center of China-made High Performance Data Computing System, Guangdong Laboratory of Artificial-Intelligence and Cyber-Economics, College of Computer Science and Software Engineering, Shenzhen University, Shenzhen, China.

A. F. Frangi is with the CISTIB Centre for Computational Imaging & Simulation Technologies in Biomedicine, School of Computing and School of Medicine, University of Leeds, Leeds, UK, and with Departments of Cardiovascular Sciences and Electrical Engineering, KU Leuven, Leuven, Belgium. He is also Pengcheng Visiting Scholar at the School of Biomedical Engineering, Shenzhen University, Shenzhen, China.

Y. Xu is with Cixi Institute of Biomedical Engineering, Ningbo Institute of Materials Technology and Engineering, Chinese Academy of Sciences, Ningbo, China.

J. Qin is with Centre for Smart Health, School of Nursing, Hong Kong Polytechnic University, HungHom, Hong Kong.

A. Elazab is with Misr Higher Institute for Commerce and Computers, Computer Science Department, Mansoura City, Egypt.

A. Elazab and B. Lei are with the National-Regional Key Technology Engineering Laboratory for Medical Ultrasound, Guangdong Key Laboratory for Biomedical Measurements and Ultrasound Imaging, School of Biomedical Engineering, Health Science Center, Shenzhen University, Shenzhen, China, 518060 (e-mail: leiby@szu.edu.cn).

explore multiple modalities, namely, GM, L1, and V1, to perform PD classification. Nevertheless, for the multimodal data, the small dataset size and large feature dimension typically cause overfitting problem and renders difficulty in model generalization [8]. Though deep learning has been extensively used in the medical image fields, it is difficult to obtain excellent generalization performance on small number of subjects [9].

Feature selection is an excellent measure to discover disease-related characteristics exploiting either supervised or unsupervised approaches [10-12]. For most supervised approaches, they can be individual single-task [13, 14] or united multi-task [15, 16]. Generally, the latter has better performance since this approach combines multiple related tasks to select common features jointly. However, there exist two main limitations in existing multi-task approaches. First, the selected features obtained from these approaches are usually linearly related to multi-task goals but they ignore the learning of the structural information intrinsic within the data. Second, these multi-task approaches need additional scores and label information to learn the model. On the other hand, unsupervised approaches focus more on learning the intrinsic data structure. Most unsupervised approaches are developed based on either filtering [17, 18] or embedding [19, 20]. The latter is superior in many aspects. However, there are three main shortcomings in existing embedding approaches. First, they calculate the similarity matrix among subjects and select features, respectively. But real data in the original high-dimensional space has noise and/or redundancy, which reduces the accuracy of the similarity matrix. Second, when calculating the neighbor graph, the similarity matrix among subjects generated by conventional approaches rarely represents a proper neighbor distribution. The optimal similarity matrix among subjects ought to have r -connected components, where r equals to the number of classes. Third, many embedding approaches do not take the similarity existing among features into consideration.

Meanwhile, most existing approaches use $\ell_{2,1}$ norm to conduct sparseness control [21, 22]. This norm cannot achieve adaptive sparseness according to different cases. Most existing studies only conducted PD classification [7, 23, 24]. Relatively few studies considered another essential task of clinical score regression [25, 26]. Since clinical score regression (e.g., depression, sleep, olfaction, and cognition scores) can assist doctors in staging and treating disease, these two tasks need to be conducted simultaneously. In addition, in most existing studies, classification and regression are performed only based on the baseline data [27, 28], while the longitudinal data (i.e., multi-time points data) are ignored. Owing to the persistent exacerbation of the disease, it is imperative to learn reliable classification and prediction models that meet multi-time points [29]. We highlight our contributions:

- 1) We propose a novel unsupervised learning method from longitudinal multimodal data for feature selection. The united embedding learning and sparse regression are exploited to adaptively learn the low-dimensional manifold structure and select the informative features.
- 2) We dynamically update the similarity matrices among subjects and features. The connected number among

subjects from the similarity matrix is equal to the number of classes, which can gain the intrinsic structural property of the data.

- 3) We conduct abundant experiments on the PPMI dataset to verify the effectiveness of the proposed approach. The results show that our algorithm effectively boosts the performance of classification and clinical score regression and surpasses other state-of-the-art approaches by taking full advantage of the longitudinal data.

The paper is organized as follows. In Section II, we discuss the most related work on feature selection. Detailed interpretation of our approach is introduced in Section III. Our results and discussions are shown in Section IV and V. Finally, several conclusions are recapitulated in Section VI.

II. RELATED WORK

In the literature, there are many supervised united multi-task and unsupervised embedding approaches. For instance, the multimodal multi-task (M3T) [30] approach based on $\ell_{2,1}$ norm learns a feature selection model to gain common relevant features of multiple tasks from every modality. The multimodal sparse learning (MMSL) [15] approach concurrently performs classification and regression prediction of PD based on a united multi-task feature selection function that considers the similarity of difference among rows and columns in response matrix. In [31], multimodal data is utilized to improve performance on the classification and regression prediction of Alzheimer's disease via relational regularization and discriminative learning. In [5], the authors perform joint learning from multiple relations and modalities to select the discriminative features for classification and prediction of PD. Multi-cluster feature selection (MCFS) [32] approach first calculates the nearest neighbor graph and then selects the discriminative features that best present the clustering information. Flexible manifold embedding (FME) [33] approach is a generalized model exploited by many unsupervised and semi-supervised embedding approaches to reduce feature dimensionality. Robust spectral feature selection (RSFS) [19] approach concurrently utilizes FME and ℓ_1 norm to robustly select the discriminative features. Joint embedding learning and sparse regression (JELSR) [34] approach conducts feature selection through embedding learning with sparse regression. These existing approaches suffer from some limitations. For instance, the supervised multi-task approaches ignore to learn the structural information intrinsic within the data. Previous unsupervised embedding approaches may also calculate an inaccurate similarity matrix due to the noise in the original feature space.

III. METHODOLOGY

A. System Overview

The overall flowchart of our approach is illustrated in Fig.1. First, we extract features from GM, L1, and V1 and concatenate them directly. We then use the proposed approach to perform feature selection. Finally, we exploit support vector classification (SVC) and support vector regression (SVR) models to conduct classification and regression prediction on longitudinal multimodal data.

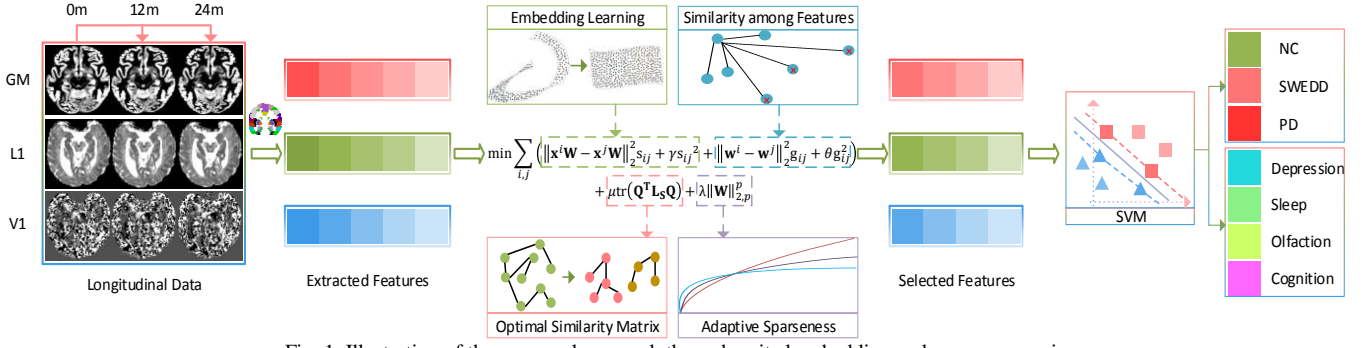


Fig. 1. Illustration of the proposed approach through united embedding and sparse regression.

B. Notation

In this study, the uppercase boldface letter \mathbf{X} indicates matrix, and the lowercase boldface letter \mathbf{x} indicates vector. Let the data matrix be $\mathbf{X} = [x_{ij}]$, \mathbf{x}^i and \mathbf{x}_j are its i -th row and j -th column, respectively. In addition, the $\ell_{2,p}$ norm of \mathbf{X} is indicated as $\|\mathbf{X}\|_{2,p} = \left(\sum_i \|\mathbf{x}^i\|_2^p\right)^{\frac{1}{p}}$. The transpose and trace operators of \mathbf{X} are indicated as \mathbf{X}^T and $tr(\mathbf{X})$, respectively.

C. Proposed Approach

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ indicate the training data of n subjects and d features and $\mathbf{S} \in \mathbb{R}^{n \times n}$ indicate the subject similarity matrix. According to our intuitive understanding, closer subjects usually have greater similarities, and thus we can calculate the similarity \mathbf{S} using the following formula:

$$\min \sum_{i,j} \left(\|\mathbf{x}^i - \mathbf{x}^j\|_2^2 s_{ij} + \mu s_{ij}^2 \right), \quad (1)$$

$$s. t. \mathbf{s}^i \mathbf{1} = 1, 0 \leq s_{ij} \leq 1$$

where \mathbf{x}^i and \mathbf{x}^j indicate the i -th and j -th subjects of \mathbf{X} , respectively. μ is a regularization parameter to avert the meaningless solution. The subject similarity matrix \mathbf{S} built with (1) rarely has a proper neighbor distribution. The optimal subject similarity matrix ought to have r -connected components, where r equals to the number of classes. Nevertheless, it is nearly impossible to obtain \mathbf{S} with (1) that satisfies the above requirements. To solve this problem, the rank of the Laplacian matrix \mathbf{L}_S of \mathbf{S} is ensure to be equal to $n - r$. In this way, there will be r -connected components in the subject similarity matrix [35]. We add this constraint on \mathbf{L}_S to (1) as:

$$\min \sum_{i,j} \left(\|\mathbf{x}^i - \mathbf{x}^j\|_2^2 s_{ij} + \mu s_{ij}^2 \right), \quad (2)$$

$$s. t. \mathbf{s}^i \mathbf{1} = 1, 0 \leq s_{ij} \leq 1, rank(\mathbf{L}_S) = n - r$$

where \mathbf{L}_S equals to $\mathbf{D} - \frac{\mathbf{S}^T + \mathbf{S}}{2}$ and the degree matrix \mathbf{D} is a diagonal matrix whose i -th element value on the diagonal equals to $\sum_j \frac{s_{ij} + s_{ji}}{2}$. Since $rank(\mathbf{L}_S) = n - r$ relies on the subject similarity matrix \mathbf{S} , it is challenging to optimize (2) directly. To solve this problem, let $\psi_i(\mathbf{L}_S)$ indicate the i -th minimum eigenvalue of \mathbf{L}_S . Because \mathbf{L}_S is positive semi-definite, we obtain $\psi_i(\mathbf{L}_S) \geq 0$. Meanwhile, it can be easy to prove that $rank(\mathbf{L}_S) = n - r$ equals $\sum_{i=1}^r \psi_i(\mathbf{L}_S) = 0$. Since it is hard to solve the derivation of $\sum_{i=1}^r \psi_i(\mathbf{L}_S)$, we refer to Ky Fan's Theorem [36], and we obtain:

$$\sum_{i=1}^r \psi_i(\mathbf{L}) = \min tr(\mathbf{Q}^T \mathbf{L}_S \mathbf{Q}), \quad (3)$$

$$s. t. \mathbf{Q} \in \mathbb{R}^{n \times r}, \mathbf{Q}^T \mathbf{Q} = \mathbf{E}$$

Further, we can rewrite (2):

$$\min \sum_{i,j} \left(\|\mathbf{x}^i - \mathbf{x}^j\|_2^2 s_{ij} + \mu s_{ij}^2 \right) + \sigma tr(\mathbf{Q}^T \mathbf{L}_S \mathbf{Q}), \quad (4)$$

$$s. t. \mathbf{s}^i \mathbf{1} = 1, 0 \leq s_{ij} \leq 1, \mathbf{Q} \in \mathbb{R}^{n \times r}, \mathbf{Q}^T \mathbf{Q} = \mathbf{E}$$

where σ is a model parameter that can be decreased or increased in each iteration to obtain the optimal \mathbf{S} when the connected components of \mathbf{S} are greater or smaller than r , respectively. \mathbf{E} is the identity matrix. In (4), the subject similarity matrix \mathbf{S} is calculated in the original multimodal feature space. However, the original high-dimensional data has noise and/or redundancy. To tackle this problem, we perform adaptive sparseness and embedding learning simultaneously, which is expressed as:

$$\min \sum_{i,j} \left(\|\mathbf{x}^i \mathbf{W} - \mathbf{x}^j \mathbf{W}\|_2^2 s_{ij} + \mu s_{ij}^2 \right) + \sigma tr(\mathbf{Q}^T \mathbf{L}_S \mathbf{Q}) + \lambda \|\mathbf{W}\|_{2,p}^p, \quad (5)$$

$$s. t. \mathbf{s}^i \mathbf{1} = 1, 0 \leq s_{ij} \leq 1, \mathbf{Q} \in \mathbb{R}^{n \times r}, \mathbf{Q}^T \mathbf{Q} = \mathbf{E}, \mathbf{W}^T \mathbf{W} = \mathbf{E}$$

where $\mathbf{W} \in \mathbb{R}^{d \times m}$ indicates the feature weight coefficient matrix of m projection dimension. λ is a weighting parameter that decreases the feature weight \mathbf{W} to obtain more sparse features as the value of λ increases. Meanwhile, we use $\ell_{2,p}$ norm to carry out sparse adaptive control for selecting the most discriminative features via different scenarios. Furthermore, the data of high-dimensional features might render the covariance matrix of \mathbf{X} singular, so we add the constraint $\mathbf{W}^T \mathbf{W} = \mathbf{E}$ to obtain the discriminative features after dimension reduction.

There exists similarity among features extracted from regions of interest (ROI). If two features, e.g., \mathbf{x}_i and \mathbf{x}_j , are similar, their weight coefficients, e.g., \mathbf{w}^i and \mathbf{w}^j , will be similar because the i -th feature \mathbf{x}_i and j -th feature \mathbf{x}_j in \mathbf{X} correspond to the i -th row \mathbf{w}^i and j -th row \mathbf{w}^j in \mathbf{W} , respectively. To exploit the relationship among features, we propose a regularization term considering the similarity among features, which is indicated as:

$$\min \sum_{i,j} \left(\|\mathbf{w}^i - \mathbf{w}^j\|_2^2 g_{ij} + \theta g_{ij}^2 \right), s. t. \mathbf{g}^i \mathbf{1} = 1, 0 \leq g_{ij} \leq 1, \quad (6)$$

where g_{ij} indicates an element in the feature similarity matrix $\mathbf{G} \in \mathbb{R}^{d \times d}$. θ is a regularization parameter to avert the meaningless solution. Finally, we add this regularization term to (5) and then we get:

$$\min \sum_{i,j} \left(\|\mathbf{x}^i \mathbf{W} - \mathbf{x}^j \mathbf{W}\|_2^2 s_{ij} + \mu s_{ij}^2 + \|\mathbf{w}^i - \mathbf{w}^j\|_2^2 g_{ij} + \theta g_{ij}^2 \right) + \sigma tr(\mathbf{Q}^T \mathbf{L}_S \mathbf{Q}) + \lambda \|\mathbf{W}\|_{2,p}^p, \quad (7)$$

$$s. t. \mathbf{s}^i \mathbf{1} = 1, 0 \leq s_{ij} \leq 1, \mathbf{g}^i \mathbf{1} = 1, 0 \leq g_{ij} \leq 1, \mathbf{Q} \in \mathbb{R}^{n \times r}, \mathbf{Q}^T \mathbf{Q} = \mathbf{E}, \mathbf{W}^T \mathbf{W} = \mathbf{E}$$

Algorithm 1: Solution to (12)

Input: $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{L}_S \in \mathbb{R}^{n \times n}$, $\mathbf{L}_G \in \mathbb{R}^{d \times d}$, λ ;
Output: $\mathbf{W} \in \mathbb{R}^{d \times m}$;
1 Initialize $t = 0$, $\mathbf{Z}(t)$ by an identity matrix;
2 **Repeat**
3 Under the current $\mathbf{Z}(t)$, the optimal solution $\mathbf{W}(t+1)$ of (12) is the m eigenvectors corresponding to the m minimum eigenvalues of $\mathbf{X}^T \mathbf{L}_S \mathbf{X} + \mathbf{L}_G + \lambda \mathbf{Z}$.
4 Update matrix $\mathbf{Z}(t+1)$ by solving (11);
7 **Until** the convergence or stop condition is satisfied.

D. Optimization

Since (7) contains $\ell_{2,p}$ norm and four variables, it is difficult to solve this problem directly. Thus, we exploit an alternative approach to tackle this problem. Meanwhile, regarding the optimization of parameter \mathbf{W} , it is necessary to use Laplacian matrices to transform (7) into trace forms. Next, we fix the similarity matrix \mathbf{S} among subjects and the similarity matrix \mathbf{G} among features and then update \mathbf{W} . Equation (7) is transformed into:

$$\min tr(\mathbf{W}^T \mathbf{X}^T \mathbf{L}_S \mathbf{X} \mathbf{W}) + tr(\mathbf{W}^T \mathbf{L}_G \mathbf{W}) + \lambda \|\mathbf{W}\|_{2,p}^p, \quad (8)$$

$s. t. \mathbf{W}^T \mathbf{W} = \mathbf{E}$

The objective function with the Lagrange multiplier for (8) is:

$$\mathcal{F}(\mathbf{W}, \Lambda) = tr(\mathbf{W}^T \mathbf{X}^T \mathbf{L}_S \mathbf{X} \mathbf{W}) + tr(\mathbf{W}^T \mathbf{L}_G \mathbf{W}) + \lambda \|\mathbf{W}\|_{2,p}^p + tr(\Lambda(\mathbf{W}^T \mathbf{W} - \mathbf{E})), \quad (9)$$

where $\Lambda \in \mathbb{R}^{m \times m}$ is the diagonal matrix indicating the Lagrange multiplier. We set the derivative of (9) on \mathbf{W} to zero and get:

$$\frac{\partial \mathcal{F}(\mathbf{W}, \Lambda)}{\partial \mathbf{W}} = \mathbf{X}^T \mathbf{L}_S \mathbf{X} \mathbf{W} + \mathbf{L}_G \mathbf{W} + \lambda \mathbf{Z} \mathbf{W} + \mathbf{W} \Lambda = 0, \quad (10)$$

where $\mathbf{Z} \in \mathbb{R}^{d \times d}$ is a diagonal matrix and its i -th diagonal element is defined as:

$$z_{ii} = \frac{p}{2(\mathbf{w}^i \mathbf{w}^{iT} + eps)^{\frac{2-p}{2}}} \quad (11)$$

where we add a floating number eps equal to 2^{-52} to the denominator of (9) since $\mathbf{w}^i \mathbf{w}^{iT}$ can be zero in theory. The solution of (10) equals to:

$$\min tr(\mathbf{W}^T \mathbf{X}^T \mathbf{L}_S \mathbf{X} \mathbf{W}) + tr(\mathbf{W}^T \mathbf{L}_G \mathbf{W}) + \lambda tr(\mathbf{W}^T \mathbf{Z} \mathbf{W}), \quad (12)$$

$s. t. \mathbf{W}^T \mathbf{W} = \mathbf{E}$

where we can exploit an iterative algorithm to solve (12) since \mathbf{Z} relies on \mathbf{W} . The details of the proposed algorithm are shown in Algorithm 1. Next, we fix the subject similarity matrix \mathbf{S} and then update \mathbf{Q} . Equation (7) is transformed into:

$$\min tr(\mathbf{Q}^T \mathbf{L}_S \mathbf{Q}), \quad (13)$$

$s. t. \mathbf{Q} \in \mathbb{R}^{n \times r}, \mathbf{Q}^T \mathbf{Q} = \mathbf{E}$

where the optimal solution \mathbf{Q} of (13) is the r eigenvectors corresponding to the r minimum eigenvalues of \mathbf{L}_S . Next, we fix \mathbf{W} and \mathbf{Q} , and then update \mathbf{S} . Equation (7) is transformed into:

$$\min \sum_{i,j} \left(\|\mathbf{x}^i \mathbf{W} - \mathbf{x}^j \mathbf{W}\|_2^2 s_{ij} + \mu s_{ij}^2 \right) + \sigma tr(\mathbf{Q}^T \mathbf{L}_S \mathbf{Q}), \quad (14)$$

$s. t. \mathbf{s}^i \mathbf{1} = 1, 0 \leq s_{ij} \leq 1,$

We can transform (14) into the following form:

$$\min \sum_{i,j} \left(\|\mathbf{x}^i \mathbf{W} - \mathbf{x}^j \mathbf{W}\|_2^2 s_{ij} + \mu s_{ij}^2 \right) + \sigma \sum_{i,j} \|\mathbf{q}^i - \mathbf{q}^j\|_2^2 s_{ij}, \quad (15)$$

$s. t. \mathbf{s}^i \mathbf{1} = 1, 0 \leq s_{ij} \leq 1$

The similarity vector of every subject is independent. Thus we can solve the above problem for the i -th subject, and we have:

$$\min \sum_j \left(\|\mathbf{x}^i \mathbf{W} - \mathbf{x}^j \mathbf{W}\|_2^2 s_{ij} + \mu s_{ij}^2 \right) + \sigma \sum_j \|\mathbf{q}^i - \mathbf{q}^j\|_2^2 s_{ij}, \quad (16)$$

$s. t. \mathbf{s}^i \mathbf{1} = 1, 0 \leq s_{ij} \leq 1$

where we indicate $\mathbf{a}^i \in \mathbb{R}^{1 \times n}$ with $a_{ij} = \|\mathbf{x}^i \mathbf{W} - \mathbf{x}^j \mathbf{W}\|_2^2 + \sigma \|\mathbf{q}^i - \mathbf{q}^j\|_2^2$. We can rewrite (16):

$$\min \left\| \mathbf{s}^i + \frac{1}{2\mu} \mathbf{a}^i \right\|_2^2, s. t. \mathbf{s}^i \mathbf{1} = 1, 0 \leq s_{ij} \leq 1 \quad (17)$$

We consider two extreme cases of μ in (1). When $\mu = 0$, it makes only a single element of \mathbf{s}^i not equal to zero. When $\mu = \infty$, it makes each element of \mathbf{s}^i equal to $\frac{1}{n}$. The value of μ controls the number of neighbors of a subject. Thus the optimal value of μ needs to make most \mathbf{s}^i have k non-zero elements, where k indicates the number of neighbors connected to \mathbf{x}^i . To achieve this goal, we consider the Lagrangian function of (17) as follows:

$$\mathcal{F}(\mathbf{s}^i, \beta, \boldsymbol{\varphi}^i, \boldsymbol{\tau}^i) = \frac{1}{2} \left\| \mathbf{s}^i + \frac{1}{2\mu} \mathbf{a}^i \right\|_2^2 + \beta (\mathbf{s}^i \mathbf{1} - 1) + \boldsymbol{\varphi}^i (\mathbf{s}^{iT} - \mathbf{1}) + \boldsymbol{\tau}^i (-\mathbf{s}^i \mathbf{1}), \quad (18)$$

where β , $\boldsymbol{\varphi}^i$, and $\boldsymbol{\tau}^i$ indicate Lagrangian multipliers. According to the Karush–Kuhn–Tucker (KKT) condition, the optimal solution of \mathbf{s}^i is:

$$s_{ij} = -\frac{1}{2\mu} a_{ij} - \beta, \quad (19)$$

Considering the convenience of expression, we suppose that $a'_{i1}, a'_{i2}, \dots, a'_{in}$ are sorted from small to large corresponding to $s'_{i1}, s'_{i2}, \dots, s'_{in}$. When the optimal \mathbf{s}^i contains only k neighbors, we know $s'_{ij} > 0$ ($j = 1, 2, \dots, k$) and $s'_{ij} = 0$ ($j = k+1, k+2, \dots, n$). Therefore, we have:

$$\begin{cases} -\frac{1}{2\mu} a'_{ij} - \beta > 0 & (j = 1, 2, \dots, k) \\ -\frac{1}{2\mu} a'_{ij} - \beta \leq 0 & (j = k+1, k+2, \dots, n) \end{cases} \quad (20)$$

According to (19) and the constraint $\mathbf{s}^i \mathbf{1} = 1$, we get:

$$\sum_{j=1}^k \left(-\frac{1}{2\mu} a'_{ij} - \beta \right) = 1 \rightarrow \beta = -\frac{1}{k} - \sum_{j=1}^k \frac{1}{2k\mu} a'_{ij}, \quad (21)$$

$s. t. j = 1, 2, \dots, k,$

where the inequality for μ is given by (20) and (21):

$$\frac{k}{2} a'_{ik} - \frac{1}{2} \sum_{j=1}^k a'_{ij} < \mu \leq \frac{k}{2} a'_{ik+1} - \frac{1}{2} \sum_{j=1}^k a'_{ij}, \quad (22)$$

To obtain a good μ that can make the most \mathbf{s}^i has k neighbors or non-zeros elements, we can calculate μ as:

$$\mu = \frac{1}{n} \sum_{i=1}^n \left(\frac{k}{2} a'_{ik+1} - \frac{1}{2} \sum_{j=1}^k a'_{ij} \right), \quad (23)$$

Finally, we fix \mathbf{W} and then update \mathbf{G} . Equation (7) is transformed into:

$$\min \sum_{i,j} \left(\|\mathbf{w}^i - \mathbf{w}^j\|_2^2 g_{ij} + \theta g_{ij}^2 \right), \quad (24)$$

$s. t. \mathbf{g}^i \mathbf{1} = 1, 0 \leq g_{ij} \leq 1$

The similarity vector of every feature is independent. Thus we can solve the above problem for the i -th feature, and we get:

$$\min \sum_j \left(\|\mathbf{w}^i - \mathbf{w}^j\|_2^2 g_{ij} + \theta g_{ij}^2 \right), \quad (25)$$

$s. t. \mathbf{g}^i \mathbf{1} = 1, 0 \leq g_{ij} \leq 1$

where we indicate $\mathbf{b}^i \in \mathbb{R}^{1 \times d}$ with $b_{ij} = \|\mathbf{w}^i - \mathbf{w}^j\|_2^2$. Therefore, we can rewrite (24) as:

Algorithm 2: Solution to (7)**Input:** $\mathbf{X} \in \mathbb{R}^{d \times n}$, r, m, k, λ ;**Output:** $\mathbf{W} \in \mathbb{R}^{d \times m}$;1 Initialize \mathbf{S} by solving (1);2 Initialize \mathbf{G} by solving the following formula:

$$\min \sum_{i,j} \left(\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 g_{ij} + \theta g_{ij}^2 \right),$$

$$s. t. \mathbf{g}^i \mathbf{1} = 1, 0 \leq g_{ij} \leq 1;$$

3 **Repeat**4 Update \mathbf{W} by Algorithm 1;5 Calculate the Laplacian matrix \mathbf{L}_S of the similarity matrix6 \mathbf{S} ;7 Update \mathbf{Q} by solving (13); Update \mathbf{S} by solving (17); Update \mathbf{G} by solving (25);8 **Until** the convergence or stop condition is satisfied.

$$\min \left\| \mathbf{g}^i + \frac{1}{2\theta} \mathbf{b}^i \right\|_2^2, s. t. \mathbf{g}^i \mathbf{1} = 1, 0 \leq g_{ij} \leq 1 \quad (26)$$

The solution of (26) is the same as solving (17). We show the details of the total optimization in Algorithm 2.

E. Convergence Analysis of Algorithm 1

To prove that algorithm 1 converges, we first need to verify the following inequality.

Theorem 1: For any positive real numbers u and v , with p a constant ($0 < p < 2$), it holds:

$$u^{\frac{p}{2}} - \frac{pu}{2v^{\frac{2-p}{2}}} \leq v^{\frac{p}{2}} - \frac{pv}{2v^{\frac{2-p}{2}}} \quad (27)$$

Proof: We move the left side of the inequality (27) to its right-hand side, and then we need to prove the following inequality:

$$\begin{aligned} \tau(u, v) &= v^{\frac{p}{2}} - \frac{pv}{2v^{\frac{2-p}{2}}} - u^{\frac{p}{2}} + \frac{pu}{2v^{\frac{2-p}{2}}} \geq 0 \\ \rightarrow \tau(u, v) &= (2-p)v + pu - 2u^{\frac{p}{2}}v^{\frac{2-p}{2}} \geq 0, s. t. u > 0, v > 0 \end{aligned} \quad (28)$$

To prove the inequality (28), we consider its Lagrangian function as:

$$\begin{aligned} \tau(u, v, \gamma_1, \gamma_2) &= (2-p)v + pu - 2u^{\frac{p}{2}}v^{\frac{2-p}{2}} - \gamma_1 u - \gamma_2 v, \quad (29) \\ \text{where } \gamma_1 \text{ and } \gamma_2 &\text{ indicate Lagrangian multipliers. According to the KKT condition, when } u \text{ equals to } v, \text{ (29) has the minimum value equal to zero. Thus, we can get:} \end{aligned}$$

$$\tau(u, v) \geq 0 \rightarrow u^{\frac{p}{2}} - \frac{pu}{2v^{\frac{2-p}{2}}} \leq v^{\frac{p}{2}} - \frac{pv}{2v^{\frac{2-p}{2}}} \quad (30)$$

which concludes the proof.

Theorem 2: The iterative updating rules in Algorithm 1 will gradually reduce the objective value of (8) until convergence.

Proof: Supposing the updated \mathbf{W} is $\widehat{\mathbf{W}}$, we can easily get:

$$\begin{aligned} &tr(\widehat{\mathbf{W}}^T (\mathbf{X}^T \mathbf{L}_S \mathbf{X} + \mathbf{L}_G) \widehat{\mathbf{W}}) + \lambda tr(\widehat{\mathbf{W}}^T \mathbf{Z} \widehat{\mathbf{W}}) \\ &\leq tr(\mathbf{W}^T (\mathbf{X}^T \mathbf{L}_S \mathbf{X} + \mathbf{L}_G) \mathbf{W}) + \lambda tr(\mathbf{W}^T \mathbf{Z} \mathbf{W}), \end{aligned} \quad (31)$$

where we add $\lambda \sum_i \frac{p \cdot eps}{2(\mathbf{w}^i \mathbf{w}^{iT} + eps)^{\frac{2-p}{2}}}$ to both sides of the inequality (31) and replace the definition of \mathbf{Z} with (11), and then we can rewrite the inequality (31):

$$tr(\widehat{\mathbf{W}}^T (\mathbf{X}^T \mathbf{L}_S \mathbf{X} + \mathbf{L}_G) \widehat{\mathbf{W}}) + \lambda \sum_i \frac{p(\widehat{\mathbf{w}}^i \widehat{\mathbf{w}}^{iT} + eps)}{2(\mathbf{w}^i \mathbf{w}^{iT} + eps)^{\frac{2-p}{2}}}$$

$$\leq tr(\mathbf{W}^T (\mathbf{X}^T \mathbf{L}_S \mathbf{X} + \mathbf{L}_G) \mathbf{W}) + \lambda \sum_i \frac{p(\mathbf{w}^i \mathbf{w}^{iT} + eps)}{2(\mathbf{w}^i \mathbf{w}^{iT} + eps)^{\frac{2-p}{2}}}, \quad (32)$$

Based on Theorem 1, we have:

$$\begin{aligned} &\lambda \sum_i \left(\widehat{\mathbf{w}}^i \widehat{\mathbf{w}}^{iT} + eps \right)^{\frac{p}{2}} - \lambda \sum_i \frac{p(\widehat{\mathbf{w}}^i \widehat{\mathbf{w}}^{iT} + eps)}{2(\mathbf{w}^i \mathbf{w}^{iT} + eps)^{\frac{2-p}{2}}} \\ &\leq \lambda \sum_i \left(\mathbf{w}^i \mathbf{w}^{iT} + eps \right)^{\frac{p}{2}} - \lambda \sum_i \frac{p(\mathbf{w}^i \mathbf{w}^{iT} + eps)}{2(\mathbf{w}^i \mathbf{w}^{iT} + eps)^{\frac{2-p}{2}}} \end{aligned} \quad (33)$$

We add inequalities (31) and (32) together, and then we can get:

$$\begin{aligned} &tr(\widehat{\mathbf{W}}^T (\mathbf{X}^T \mathbf{L}_S \mathbf{X} + \mathbf{L}_G) \widehat{\mathbf{W}}) + \lambda \sum_i \left(\widehat{\mathbf{w}}^i \widehat{\mathbf{w}}^{iT} + eps \right)^{\frac{p}{2}} \\ &\leq tr(\mathbf{W}^T (\mathbf{X}^T \mathbf{L}_S \mathbf{X} + \mathbf{L}_G) \mathbf{W}) + \sum_i \left(\mathbf{w}^i \mathbf{w}^{iT} + eps \right)^{\frac{p}{2}} \end{aligned} \quad (34)$$

which concludes the proof.

IV. EXPERIMENT

A. Data Acquisition

The multimodal data used in this paper is obtained from the PPMI database. It is the first comprehensive, large-scale, and international database to study PD. In this study, we employ MRI and DTI data acquired by the Siemens MAGNETOM Trio 3.0 T MRI scanner. We select MRI data using these parameters: field strength= 3 tesla, flip angle= 9°, slice thickness =1 mm, TR = 2300 ms, TE= 2.98 ms, pulse sequence = GR/IR, and acquisition plane = SAGITTAL. For DTI data, the data acquisition parameters are: field strength = 3 tesla, flip angle = 90°, slice thickness = 2 mm, gradient directions = 64, TR = 600-1000 ms, and TE = 88 ms, pulse sequence = EP.

B. Subjects

In this paper, we collect baseline MRI and DTI data acquired from 238 subjects, including 62 normal control (NC) subjects, 142 PD subjects, and 34 SWEDD subjects. We collect 12-month data acquired from 186 subjects, including 54 NC subjects, 123 PD subjects, and 9 SWEDD subjects. We collect 24-month data acquired from 127 subjects, including 7 NC subjects, 98 PD subjects, and 22 SWEDD subjects. Geriatric Depression Scale (GDS), Epworth Sleepiness Scale (ESS), University of Pennsylvania Smell Identification Test (UPSIT), and Montreal Cognitive Assessment (MoCA) are used to estimate the depression, sleep, olfaction, and cognition scores, respectively. The GDS score is estimated according to a yes/no survey's answer. The GDS scores of normal condition, mild, moderate, and severe depression are between 0 and 4, 5 and 7, 8 and 11, 12 and 15, respectively. The ESS score is estimated according to a weighted sum of responses to several questions. The ESS scores of normal and sleepy subjects are between 0 and 9, 10 and 24, respectively. The UPSIT score is between 0 and 40. A low UPSIT score means that the subject has lost a lot of smell sense. The MoCA score is between 0 and 30. A low MoCA score means that the subject has lost a lot of cognitive ability. Table 1 summarizes the clinical information of the subjects.

TABLE I
CLINICAL DETAILS OF ALL SUBJECTS IN LONGITUDINAL TIME POINT

Time	Information	NC	PD	SWEDD
Baseline	Number	62	142	34
	GDS scores	5.1±1.2	5.3±1.5	5.6±1.3
	ESS scores	6.5±3.9	5.9±3.4	8.7±4.2
	UPSIT scores	33.4±4.7	22.5±8.5	29.9±8.4
	MoCA scores	28.2±1.1	27.5±2.1	26.9±2.9
12m	Number	54	123	9
	GDS scores	4.9±1.4	5.3±1.3	5.6±1.0
	ESS scores	6.1±3.7	6.6±4.3	7.2±3.7
	UPSIT scores	-	-	-
	MoCA scores	27.5±1.9	26.8±2.9	26.4±2.9
24m	Number	7	98	22
	GDS scores	4.9±0.4	5.7±1.6	5.4±1.5
	ESS scores	7.1±3.6	7.9±4.4	7.2±3.6
	UPSIT scores	-	-	-
	MoCA scores	28.3±1.0	26.6±2.9	26.0±2.7

C. Data Preprocessing

For MRI data, we conduct anterior commissure-posterior commissure reorientation exploiting center-of-mass method [37] for all images and use statistical parametric mapping (SPM8) (<http://www.fil.ion.ucl.ac.uk/spm>) tool to perform the preprocessing procedures based on a well-accepted pipeline. First, we correct the head movement and geometric distortion and then use the graph-cut method [38] to conduct skull-stripping. Then, all MRI images are registered with the international consortium for brain mapping template that provides coordinates of the relevant anatomical labels. After that, we segment the corresponding anatomical regions into GM, white matter, and cerebrospinal fluid (CSF). Meanwhile, these images are resampled to an isotropic resolution of 1.5mm. We spatially smooth the surface of these images using a 60-mm full width at half-maximum Gaussian kernel. The purpose of smoothing is to suppress the interference of noise. Finally, we use a toolbox for data processing and analysis for brain imaging (<http://rfmri.org/dpabi>) to register automated anatomical labeling (AAL) atlas [39] with GM and extract 116-dimensional features from GM exploiting the registered AAL atlas.

For DTI data, each subject contains 65 original format images where the b0 image does not activate the diffusion gradient, while the other 64 images have different gradient directions. First, we use the FMRIB Software Library (FSL) [40] to correct the b0 image distortion. Second, we use *bet* command of FSL tool to generate a mask image corresponding to the corrected b0 image. Third, we use *dcm2nii* tool (https://www.nitrc.org/frs/?group_id=152) to convert the 65 images into a 4D image and generate a b-vector file and a b-value file indicating each gradient direction and its scalar value, respectively. Fourth, we use *eddy_correct* command of FSL tool to correct the eddy current distortion on the 4D image. Fifth, we import the b-value file, the b-vector file, the mask image, and the corrected 4D image to *dtifit* command of FSL tool to calculate L1 and V1 images. Finally, we use AAL atlas to calculate the mean tissue density of each region of L1 and V1 and then obtain their 116-dimensional features, respectively.

D. Experimental Setting

We use the 10-fold cross-validation approach to verify the proposed approach in baseline multimodal data (GM, L1, and V1). Specifically, we randomly separate the baseline dataset in ten groups, where one group is used for testing, and the rest is

used for training. We duplicate this process ten times to avert the probable bias during data partition. The final result is calculated by averaging the above results. We perform experiments on three binary classifications (i.e., NC vs. PD, NC vs. SWEDD, and PD vs. SWEDD) and prediction of four scores (i.e., depression, sleep, olfaction, and cognition scores) in baseline multimodal data. Due to the existence of data loss problem on longitudinal time points, we use the 5-fold cross-validation approach to verify the proposed approach on the 12-month and 24-month data. In addition, to enhance the generalization ability, we use baseline data as part of the training data to help 12-month data to learn the classification and regression prediction models and exploit baseline and 12-month data as part of the training data to assist 24-month data in determining the classification and regression prediction models. We also conduct three binary classification experiments on the 12-month and 24-month data. Due to the lack of olfaction scores on the 12-month and 24-month data, we only predict depression, sleep, and cognition scores. We determine the optimal SVC/SVR parameters of the support vector machine from $S_C \in \{2^{-10}, \dots, 2^{10}\}$, and $S_G \in \{2^{-5}, \dots, 2^5\}$, by performing grid search on the hyper-parameters of our objective function with the spaces of $\lambda \in \{10^0, \dots, 10^8\}$, $m \in \{100, 120, \dots, 240\}$, $k \in \{1, 2, \dots, 10\}$, and $p \in \{0.1, 0.4, \dots, 1.9\}$. Other parameters of the objective function can be adaptively determined during the model optimization.

E. Algorithm Comparison

The proposed approach is compared with state-of-the-art approaches including (1) principal component analysis (PCA) [18], which is added into the MATLAB software as an unsupervised dimensionality reduction method; (2) Laplacian score (Lscore) [17] (<http://www.cad.zju.edu.cn/home/dengcai/>), which is an unsupervised feature selection approach (<http://www.cad.zju.edu.cn/home/dengcai/>); the core idea of Lscore is to estimate the features based on their locality preserving ability; (3) RSFS [19] (<https://github.com/LeiShiCS/RSFS>), which concurrently exploits FME and ℓ_1 norm to robustly select the discriminative features; (4) the MMSL approach obtained from Lei *et al.* [15], which concurrently conducts classification and regression prediction of PD based on a united multi-task feature selection function that considers the similarity of difference among rows and columns in response matrix; (5) the joint multi-task learning (JMTL) approach from Lei *et al.* [5], which performs classification and prediction of PD based on a united multi-task feature selection function that explores multiple relationships in the response matrix, and (6) the M3T [30] approach based on $\ell_{2,1}$ norm, which learns a feature selection model to gain common relevant features of multiple tasks from every modality; the M3T approach is a particular case of MMSL approach when its two regularization terms set to zero.

F. Model Training

For PCA, we first learn the principal component coefficient matrix, and then we multiply the original data and the coefficient matrix to conduct the feature dimension reduction. Meanwhile, for a fair comparison, other approaches select features: we first calculate the ℓ_2 norm of each row of feature weight matrix \mathbf{W} to obtain the column vector \mathbf{b} and then get its

average value c . Based on the empirical results, we select features corresponding to element values of b greater than or equal to $0.2 \times c$. We import the selected features into the support vector machine (<https://www.csie.ntu.edu.tw/~cjlin/libsvm/>) to learn SVC models for the classification and SVR models for the clinical score prediction.

G. Evaluation Criteria

To estimate the performance of the competing approaches, we use the quantitative measurements of accuracy (ACC), sensitivity (SEN), specificity (SPEC), precision (PREC), unweighted average recall (UAR), F1-score, and area under the receiver operating characteristic (ROC) curve (AUC) in classification tasks, and also the quantitative measurements of Pearson's correlation coefficient (CC), root mean squared error (RMSE), and mean absolute error (MAE) in regression tasks.

H. Classification Performance

Tables II, III, and IV show the classification performance of the competing approaches on longitudinal multimodal data. Meanwhile, the ROC curves of the related approaches are compared on three binary classification tasks in Fig. 2. According to these results, we can obtain the following findings.

First, unsupervised embedding approaches, such as RSFS and the proposed approach, have better performance than filter approaches such as Lscore and PCA. For example, regarding

NC vs. SWEDD classification on the baseline data, RSFS reaches ACC of 85.56%, F1-score of 90.23%, and AUC of 0.75, while our approach reaches ACC of 89.78%, F1-score of 93.28%, and AUC of 0.84. However, PCA achieves ACC only of 80.22%, F1-score of 86.07%, and AUC of 0.75 and Lscore merely achieves ACC of 83.67%, F1-score of 88.68%, and AUC of 0.75. We see that PCA performs best with 24-month data in NC vs. SWEDD and obtains ACC of 82.48%, F1-score of 76%, and AUC of 0.95, and the other approaches perform relatively worse. The main reason may be that the sample size of 24-month data is too small in NC vs. SWEDD. Compared with Lscore, RSFS, M3T, MMSL, JMTL, and the proposed approaches that learn feature weight to select the discriminative features, PCA exploits feature coding to achieve feature dimensionality reduction, and thus it would be more efficient to perform better under limited data.

Second, unsupervised approaches are harder than supervised approaches because the label information is missing. However, the proposed approach has better classification performance than M3T, MMSL, and JMTL. For example, our approach achieves accuracies higher than M3T, MMSL, and JMTL on the baseline data, i.e., our approach and the other three approaches achieve accuracies of 83.33% vs. 78.90% vs. 81.33% vs. 81.81% for NC and PD, 89.78% vs. 82.44% vs. 87.44% vs. 88.67 for NC and SWEDD, and 88.69% vs. 85.85% vs. 87.52%

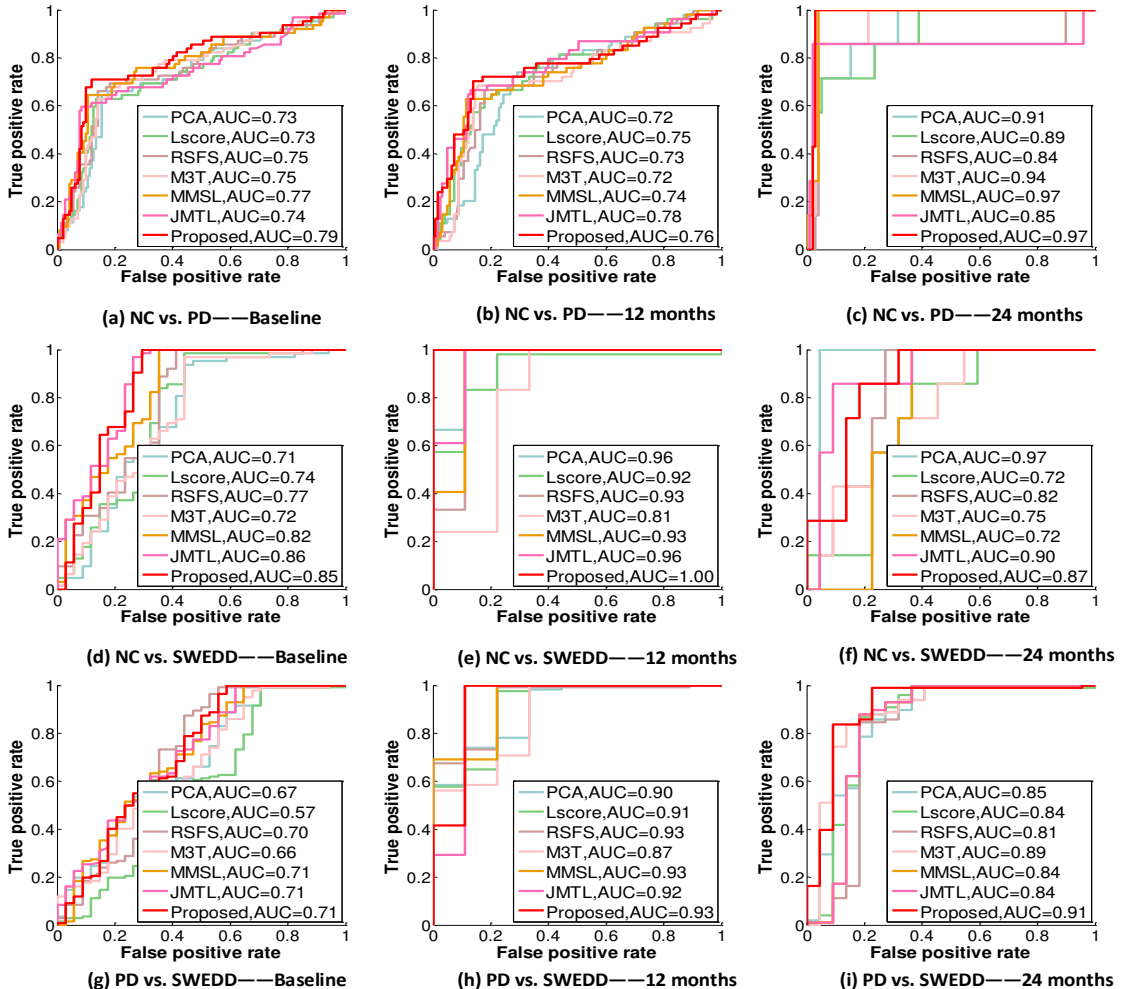


Fig. 2. Comparison of ROC curves for the competing approaches on longitudinal multimodal data.

Proposed		99.26±1.66	100.00±0.00	90.00±22.36	99.23±1.72	95.00±11.18	99.61±0.88	0.92±0.18
24m	PCA	92.49±5.46	100.00±0.00	60.00±29.37	91.81±5.78	80.00±14.68	95.65±3.11	0.88±0.11
	Lscore	92.49±5.46	98.95±2.35	65.00±23.45	92.53±5.12	81.97±12.20	95.58±3.20	0.83±0.18
	RSFS	93.32±5.61	100.00±0.00	64.00±30.70	92.71±5.88	82.00±15.35	96.14±3.17	0.82±0.24
	M3T	92.49±5.46	100.00±0.00	60.00±29.37	91.81±5.78	80.00±14.68	95.65±3.11	0.88±0.11
	MMSL	93.36±4.64	100.00±0.00	65.00±23.45	92.63±5.06	82.50±11.73	96.12±2.70	0.82±0.20
	JMTL	93.36±4.64	100.00±0.00	65.00±23.45	92.63±5.06	82.50±11.73	96.12±2.70	0.82±0.20
	Proposed	95.03±3.40	98.95±2.35	77.00±22.80	95.36±4.55	87.97±10.78	97.04±1.97	0.89±0.10

	PCA	Lscore	RSFS	M3T	MMSL	JMTL	Proposed	
0m	NC	20 22	123 19	121 21	122 20	126 16	130 12	128 14
	PD	22 40	23 39	21 41	23 39	22 40	25 37	20 42
12m	NC	93 30	97 26	101 22	105 18	109 14	108 15	106 17
	PD	19 35	19 35	18 36	18 36	20 34	19 35	16 38
24m	NC	94 4	93 5	93 5	94 4	94 4	96 2	95 3
	PD	2 5	2 5	1 6	1 6	0 7	1 6	0 7
0m	NC	18 16	19 15	20 14	19 15	22 12	23 11	24 10
	SWEDD	3 59	1 61	0 62	2 60	0 62	0 62	0 62
12m	NC	8 1	7 2	8 1	6 3	8 1	8 1	9 0
	SWEDD	1 53	1 53	0 54	0 54	0 54	0 54	0 54
24m	NC	17 5	7 15	14 8	10 12	14 8	14 8	12 10
	SWEDD	0 7	0 7	0 7	0 7	0 7	0 7	0 7
0m	PD	10 24	10 24	14 20	10 24	12 22	13 21	14 20
	SWEDD	1 141	1 141	1 141	1 141	0 142	0 142	0 142
12m	PD	5 4	6 3	7 2	6 3	7 2	7 2	8 1
	SWEDD	0 123	0 123	1 122	0 123	0 123	0 123	0 123
24m	PD	13 9	14 8	14 8	13 9	14 8	14 8	17 5
	SWEDD	0 98	1 97	0 98	0 98	0 98	0 98	1 97

Fig. 3. Confusion matrices for the competing approaches on longitudinal data.

vs. 88.07% for PD vs. SWEDD. The results demonstrate that the proposed approach can better learn the structural property intrinsic in data.

Third, on the baseline data, specificity values of the competing approaches are low for discriminating PD and SWEDD. For example, PCA, Lscore, RSFS, M3T, MMSL, JMTL, and the proposed approaches have specificity values of 28.33%, 28.33%, 39.17%, 29.17%, 34.17%, 37.50%, and 40.83%. The reasons may be that it is difficult to distinguish SWEDD patients from PD since both PD and SWEDD have asymmetric rest tremor, and the size of data is limited. When we exploit baseline data to help 12-month data to learn the classification model and exploit baseline and 12-month data to assist 24-month data in learning the classification and model, the classification performance improves. For instance, on the 12-month data, our approach achieves ACC of 99.26%, SEN of 100%, SPEC of 90%, PREC of 99.23%, UAR of 95%, F1-score of 99.61%, and AUC of 0.92. On the 24-month data, the proposed approach has ACC of 95.03%, SEN of 98.95%, SPEC of 77%, PREC of 95.36%, UAR of 87.97%, F1-score of 97.04%, and AUC of 0.89.

Finally, the proposed approach achieves the best performance on most longitudinal time points. Taking NC vs. PD as an example, on the baseline data, our approach has ACC of 83.33%, SEN of 68.33%, SPEC of 90.10%, PREC of 76.52%, UAR of 79.21%, F1-score of 71.17%, and AUC of 0.76. On the

12-month data, the proposed approach has ACC of 81.44%, SEN of 70.18%, SPEC of 86.33%, PREC of 71.89%, UAR of 78.26%, F1-score of 70.30%, and AUC of 0.76. On the 24-month data, our approach has ACC of 97.19%, SEN of 100%, SPEC of 97%, PREC of 80%, UAR of 98.5%, F1-score of 86%, and AUC of 0.98. Meanwhile, we see that our approach obtains the highest UAR overall on the longitudinal data in Tables II, III, and IV, also demonstrating the superiority of our approach. Further, we obtain the overall/average confusion matrices through linearly connecting the predicted labels and the actual labels of the test data in the entire cross-validation process. Fig.3 shows the confusion matrices of the competing approaches on longitudinal data. We can observe that our approach obtains the optimum performance.

I. Regression Performance

Tables V, VI, and VII summarize the regression performances of the competing approaches on longitudinal multimodal data. According to CC results, our approach and JMTL approach obtain the optimum regression performance overall.

In the regression performance for NC vs. PD, our approach achieves the optimal performance for olfaction and cognition scores on the baseline data. The corresponding CC, RMSE, and MAE are 0.553, 8.243, and 6.820 for olfaction score, and 0.608, 3.355, and 2.744 for cognition score. JMTL shows the optimal performance for depression score on the baseline data. The corresponding CC, RMSE, and MAE are 0.603, 4.320, and 3.466 for depression score. MMSL shows the optimal performance for sleep score on the baseline data. The corresponding CC, RMSE, and MAE are 0.569, 5.743, and 4.683. On the 12-month data, our approach achieves the optimal performance for depression and sleep scores. The corresponding CC, RMSE, MAE are 0.543, 2.083, and 1.643 for depression score, 0.514, 6.019, and 4.800 for sleep score. JMTL shows the optimal performance for cognition score. The corresponding CC, RMSE, and MAE are 0.590, 2.681, and 2.006 for cognition score. On the 24-month data, our approach has the optimal performance for depression and cognition scores. The corresponding CC, RMSE, and MAE are 0.655, 1.546, and 1.064 for depression score and 0.726, 2.113, and 1.641 for cognition score. JMTL has the optimal performance for sleep score on the 24-month data. The corresponding CC, RMSE, and MAE are 0.587, 5.425, and 4.293.

In the regression performance for NC vs. SWEDD, our approach obtains the optimal performance for sleep, olfaction, and cognition scores on the baseline data. The corresponding results are CC of 0.763, RMSE of 5.147, and MAE of 4.427 for sleep score, CC of 0.775, RMSE of 5.752, and MAE of 4.374 for olfaction score, and CC of 0.790, RMSE of 3.632, and MAE of 3.062 for cognition score. JMTL obtains the optimal performance for depression score on the baseline data. The corresponding results are CC of 0.777, RMSE of 2.473, and MAE

	Proposed	0.561	2.123	1.719	0.535	4.862	3.885	--	--	--	0.580	3.716	2.811
24m	PCA	0.511	1.452	1.064	0.441	4.814	3.750	--	--	--	0.686	3.678	3.003
	Lscore	0.491	1.563	1.123	0.392	5.568	4.382	--	--	--	0.673	3.082	2.503
	RSFS	0.576	2.713	2.193	0.501	4.420	3.373	--	--	--	0.726	3.119	2.520
	M3T	0.533	3.147	2.554	0.442	4.781	3.936	--	--	--	0.670	3.086	2.503
	MMSL	0.558	2.158	1.644	0.495	5.080	3.962	--	--	--	0.703	3.113	2.521
	JMTL	0.558	2.158	1.644	0.496	6.480	5.120	--	--	--	0.698	3.400	2.699
	Proposed	0.636	1.335	0.982	0.540	4.688	3.535	--	--	--	0.726	2.995	2.491

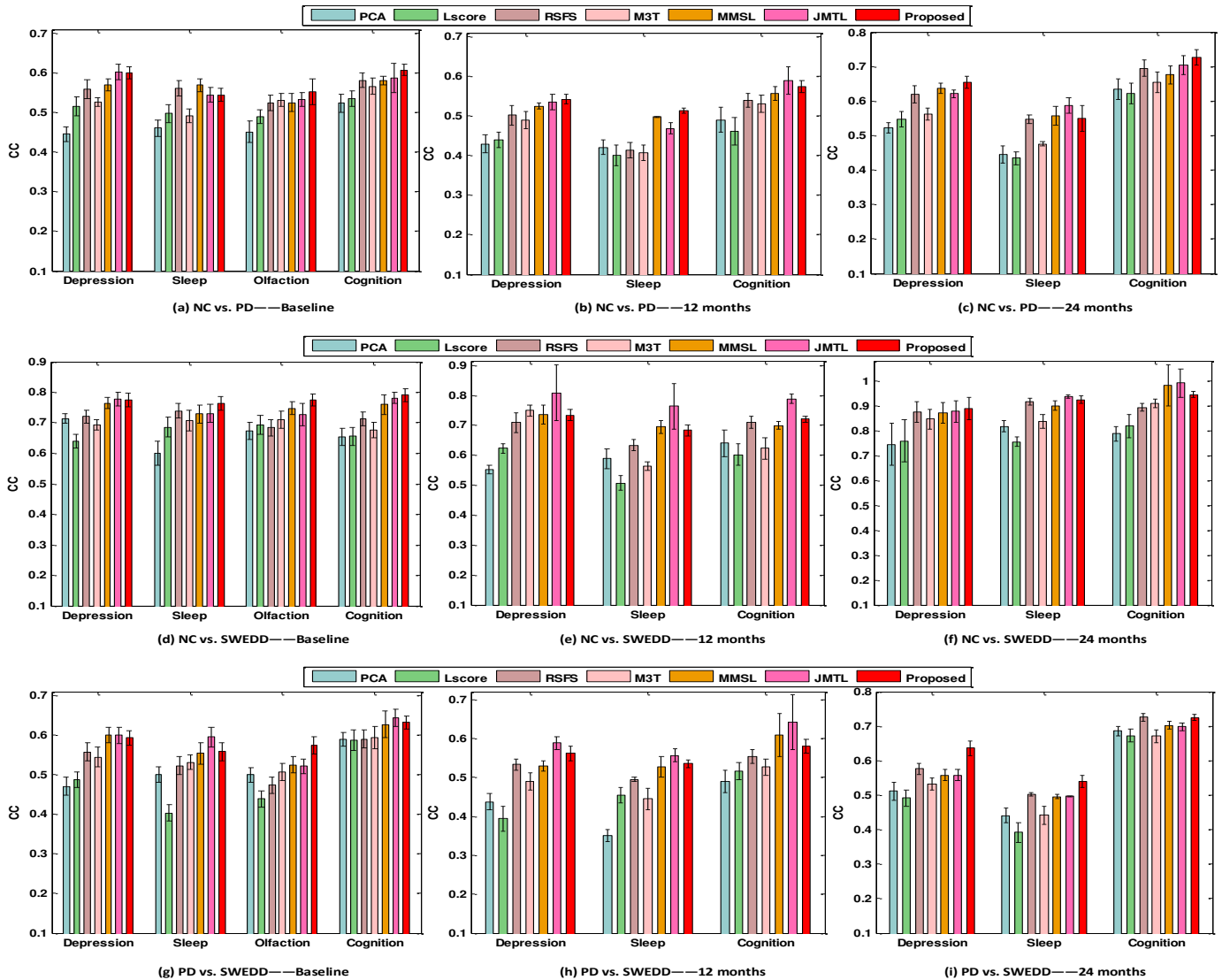


Fig. 4. Regression performance for the competing approaches on longitudinal multimodal data.

of 1.877 for depression score. On the 12-month data, JMTL shows the optimal performance for depression, sleep, and cognition scores. The corresponding results are CC of 0.809, RMSE of 1.289, and MAE of 0.887 for depression, CC of 0.763, RMSE of 3.828, and MAE of 3.122 for sleep score, and CC of 0.788, RMSE of 3.605, and MAE of 2.871 for cognition score. On the 24-month data, our approach achieves the optimal performance for depression score. The corresponding results are CC of 0.888, RMSE of 1.542, and MAE of 1.185 for depression score. JMTL shows the optimal performance for sleep and cognition scores. The corresponding results are CC of 0.938, RMSE of 7.763, and MAE of 6.491 for sleep score and CC of 0.990, RMSE of 3.128, and MAE of 2.446 for cognition score.

In the regression performance for PD vs. SWEDD, our approach shows the optimal performance for olfaction score in baseline data. The corresponding CC, RMSE, and MAE are 0.575, 8.277, and 6.945. MMSL shows the optimal performance for depression score. The corresponding CC, RMSE, and MAE are 0.601, 4.237, and 3.344. JMTL shows the optimal performance for sleep and cognition scores. The corresponding CC, RMSE, and MAE are 0.595, 5.431, and 4.323 for sleep score and 0.644, 2.890, and 2.216 for cognition score. On the 12-month data, JMTL shows the optimal performance for depression, sleep, and cognition scores. The corresponding CC, RMSE, and MAE are 0.588, 2.908, and 2.384 for depression score, 0.556, 4.546, and 3.652 for sleep score, and 0.641, 3.266,

and 2.501 for cognition score. On the 24-month data, our approach shows the optimal performance for depression, sleep, and cognition scores. The corresponding CC, RMSE, and MAE are 0.636, 1.335, and 0.982 for depression score, 0.540, 4.688, and 3.535 for sleep score, and 0.726, 2.995, and 2.491 for cognition score. Fig. 4 also illustrates regression performance for the competing approaches on longitudinal data, and our approach has good performance.

J. Results Summary

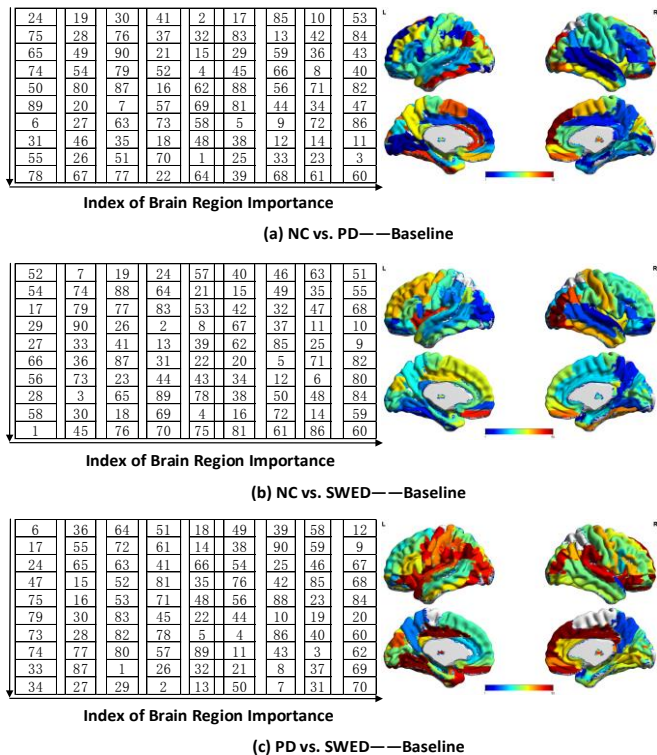


Fig. 5. Index of brain region importance and corresponding visual representation on the baseline data (Note that the importance of brain regions is sorted from top to bottom and left to right in the left table).

According to the experimental results, we observe that our approach has excellent classification performance and regression performance in NC vs. PD, NC vs. SWEDD, and PD vs. SWEDD. The proposed approach is far superior to all other approaches for the classification tasks on longitudinal data. Meanwhile, we sum up the weights of the same brain regions from three modalities to obtain the discriminative brain region index. Fig. 5. shows the index of brain region importance of baseline time point and its corresponding visual representation exploiting the BrainNet Viewer tool [41]. For NC vs. PD on the baseline data, the top ten brain regions are: right superior frontal gyrus, left pallidum, left angular gyrus, right putamen, right superior occipital gyrus, left inferior temporal gyrus, right superior frontal gyrus (orbital part), left anterior cingulate and paracingulate gyri, left fusiform gyrus, and right thalamus. For NC vs. SWEDD, the top ten brain regions are right middle occipital gyrus, right inferior occipital gyrus, left Rolandic operculum, left insula, left gyrus rectus, right angular gyrus, right fusiform gyrus, right gyrus rectus, right postcentral gyrus, and left precentral gyrus. For PD vs. SWEDD, the top ten brain regions are right superior frontal gyrus (orbital part), left Rolandic operculum, right superior frontal gyrus, left lingual

gyrus, left pallidum, left Heschl gyrus, bilateral putamen, and bilateral median cingulate and paracingulate gyri. It is noteworthy these brain regions follow the previous PD studies such as superior frontal gyrus, pallidum, angular, and putamen in [42], superior frontal gyrus (orbital part), inferior temporal gyrus, and fusiform gyrus in [43]. The top regions strongly correlate with PD, which may be the potential factors causing the disease.

To further discover the connection relationships between the top ten brain regions and other brain regions, we use the feature weights obtained by ten-fold cross-validation on our approach to calculate the Pearson correlation coefficient matrix. We then exploit the matrix to obtain five other brain regions strongly connected to each top brain region. In Fig 6, we exploit the *circularGraph* function of MATLAB software to generate the correlation maps. In NC vs. PD, we can see that some other brain regions strongly connected to multiple top brain regions follow the previous PD studies such as Heschl gyrus and superior parietal gyrus [44, 45]. We can also see that the top ten brain regions are not only strongly related to each other but also connected to other brain regions in NC vs. PD and PD vs. SWEDD. However, in [5], most of the top brain regions are only related to each other. Compared with the results in [5], our correlation maps can show richer information and thus are more suitable as a complement for PD diagnosis. The possible reason is that the multi-modal data (i.e., GM, L1, and V1) used in this paper is more able to present brain change information more effectively. For example, GM of MRI is widely used to obtain information about changes in nerve cells. L1 and V1 of DTI indicate the largest diffusion coefficient and its direction vector, respectively. Therefore, L1 and V1 may be more sensitive to neurodegeneration in the brain. Meanwhile, many current approaches [46, 47] use resting-state functional MRI data to draw correlation maps, which can effectively show the functional connection differences between different groups or stages. In this paper, we use MRI and DTI data to draw correlation maps, which can effectively show the connection relationships between the top brain regions and other related brain regions at the structural level, which is complementary to these functional studies for PD diagnosis. In addition, we can also see that the brain region importance and the correlation maps obtained by our approach have many differences from those obtained in [5]. The main reason is that we use data from different modalities in the two works. Specifically, we use GM, L1, and V1 data in this paper while using GM, CSF of MRI, and mean diffusion (MD) coefficient of DTI in [5].

Finally, Fig. 7 further illustrates the top twenty brain regions connectivity on longitudinal data, where the blue, green, and cyan lines indicate baseline and 12-month, and 24-month data, respectively. We can see that the condition worsens over time, and the brain lesion regions also change, which can verify that the proposed approach is effective. In NC vs. PD, we can also see that some top brain regions present at three-time points follow the previous PD studies such as angular gyrus, superior occipital gyrus, and inferior temporal gyrus [48, 49]. These brain regions are also related to brain cognition. At present, most studies mainly focus on PD, and there are few studies on SWEDD [5]. In Fig. 6 and Fig. 7, we also show the correlation and connection maps of NC vs. SWEDD and PD vs. SWEDD, which is helpful to guide doctors to further divide the disease.

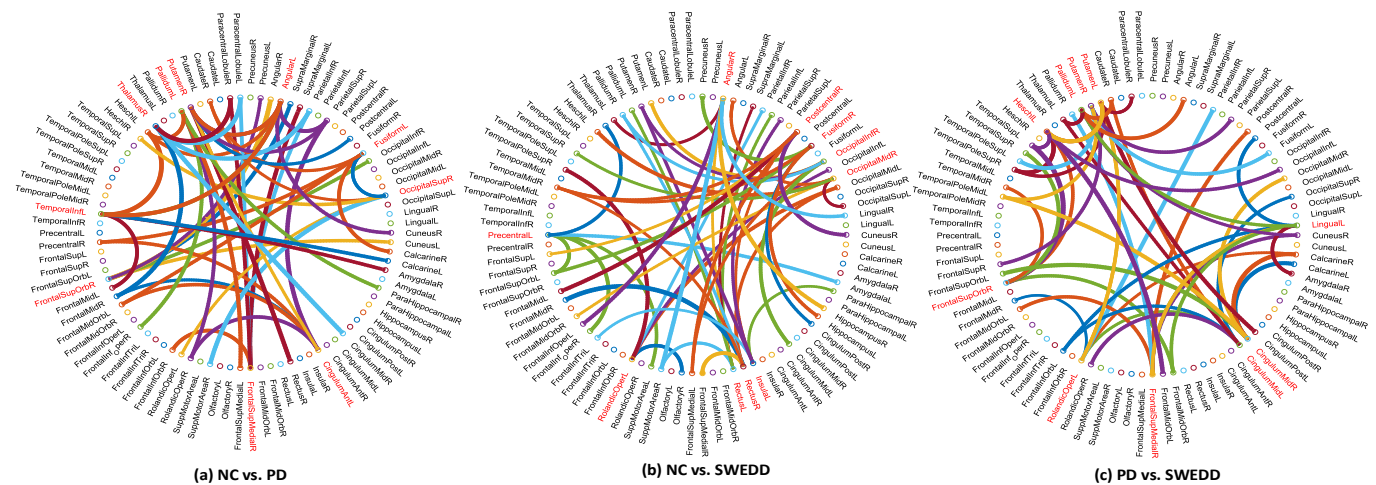


Fig. 6. Correlation maps between top ten brain regions with red fonts and other brain regions with black fonts on the baseline data.

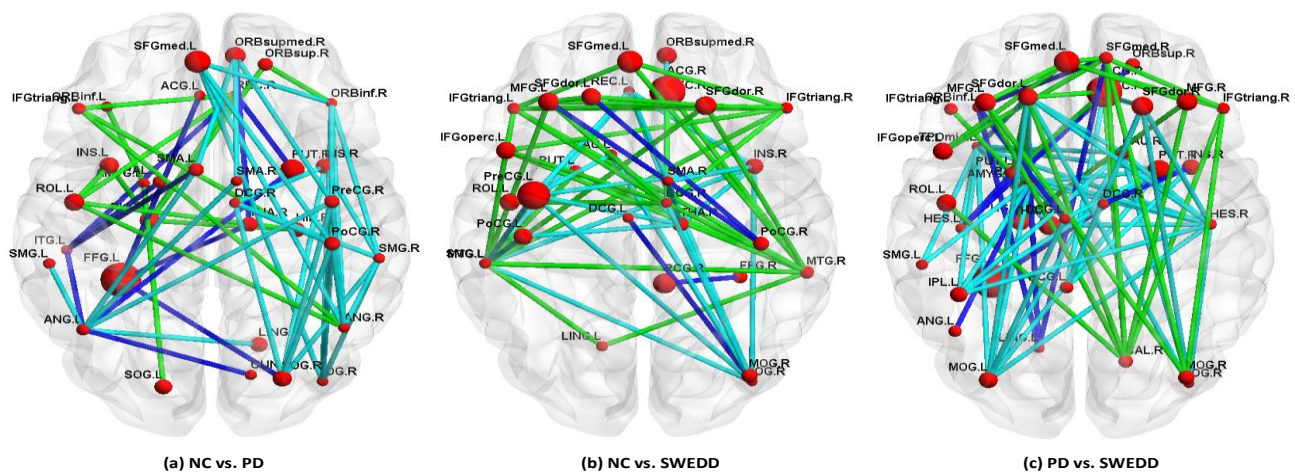


Fig. 7. Top twenty brain regions connectivity on longitudinal data.

V. DISCUSSIONS

Our approach has achieved promising performance, but there are still several limitations. First, in the data preprocessing, MRI images are registered to a standard template for tissue segmentation, which may remove some pathological changes of PD. We can use deep learning techniques to segment neuroimaging data to eliminate this effect [50-52]. Second, in this paper, we conduct the three binary classification tasks instead of a multi-class classifier for disease classification, which is consistent with the doctor's clinical diagnosis. The doctor first diagnoses whether the subject is sick and then further tests the condition of the subject in detail. Many previous studies have learned multiple binary classifiers instead of a multi-classifier for disease classification [7]. Third, our approach does not consider the relationship among longitudinal data. We are considering to add the regularization term between longitudinal data to enhance the generalization ability of the proposed approach [53]. Fourth, our work only uses MRI and DTI data to perform united classification and regression. Since some genes are related to PD [54], we can combine gene and neuroimaging data for improving the classification and regression performance. Finally, we do not analyze the full extent DTI data in this paper. For example, we can use DTI data to generate fiber bundle imaging of top brain regions and quantify white matter

fiber differences between NC and PD [55], which may clinically contribute to the classification and prediction of PD.

VI. CONCLUSIONS

In this paper, we propose a novel adaptive unsupervised feature selection approach through embedding learning using longitudinal multimodal data for the united classification and clinical score prediction of PD. Specifically, the proposed approach concurrently performs adaptive embedding learning and sparse regression; the similarity matrices and discriminative features thus can be determined adaptively. Meanwhile, we dynamically update the similarity matrices among subjects and features and have the connected number of the similarity matrix among subjects equal to the number of classes to gain the intrinsic structural property of the data. An effective iterative optimization algorithm is proposed to solve this problem. By constructing the united embedding and sparse regression framework, our approach can find the most disease-related biomarkers, which is helpful for PD monitoring. We perform abundant experiments on the PPMI dataset to verify the validity of the proposed approach. We use longitudinal data to boost the performance of regression and classification effectively. The proposed approach is shown to surpass other state-of-the-art methods.

REFERENCES

- [1] L. V. Kalia and A. E. Lang, "Parkinson's disease," *The Lancet*, vol. 386, pp. 896-912, Aug 2015.
- [2] R. B. Postuma, D. Berg, M. Stern, W. Poewe, C. W. Olanow, W. Oertel, *et al.*, "MDS clinical diagnostic criteria for Parkinson's disease," *Mov Disord*, vol. 30, pp. 1591-601, Oct 2015.
- [3] J. Lohtharius and P. Brundin, "Pathogenesis of parkinson's disease: dopamine, vesicles and [alpha]-synuclein," *Nat Rev Neurosci*, vol. 3, pp. 932-942, Jan 2002.
- [4] B. S. Connolly and A. E. Lang, "Pharmacological Treatment of Parkinson Disease: A Review," *JAMA*, vol. 311, pp. 1670-1683, 2014.
- [5] H. Lei, Z. Huang, F. Zhou, A. Elazab, E. Tan, H. Li, *et al.*, "Parkinson's Disease Diagnosis via Joint Learning From Multiple Modalities and Relations," *IEEE Journal of Biomedical and Health Informatics*, vol. 23, pp. 1437-1449, 2019.
- [6] Y. Li, J. Liu, Z. Tang, and B. Lei, "Deep Spatial-Temporal Feature Fusion From Adaptive Dynamic Functional Connectivity for MCI Identification," *IEEE Transactions on Medical Imaging*, vol. 39, pp. 2818-2830, 2020.
- [7] E. Adeli, K. Thung, L. An, G. Wu, F. Shi, T. Wang, *et al.*, "Semi-Supervised Discriminative Classification Robust to Sample-Outliers and Feature-Noises," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, pp. 515-522, 2019.
- [8] H. I. Suk, S. W. Lee, and D. Shen, "Subclass-based multi-task learning for Alzheimer's disease diagnosis," *Front Aging Neurosci*, vol. 6, p. 168, 2014.
- [9] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, *et al.*, "A survey on deep learning in medical image analysis," *Medical Image Analysis*, vol. 42, pp. 60-88, 2017.
- [10] F. Nie, W. Zhu, and X. Li, "Unsupervised feature selection with structured graph optimization," in *the Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, 2016, pp. 1302-1308.
- [11] X. Zhu, H.-I. Suk, S.-W. Lee, and D. Shen, "Subspace Regularized Sparse Multi-Task Learning for Multi-Class Neurodegenerative Disease Identification," *IEEE transactions on biomedical engineering*, vol. 63, pp. 607-618, 2016.
- [12] H. Lei, Z. Huang, A. Elazab, H. Li, and B. Lei, "Longitudinal and Multi-modal Data Learning via Joint Embedding and Sparse Regression for Parkinson's Disease Diagnosis," in *Machine Learning in Medical Imaging*, 2018, pp. 310-318.
- [13] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *Journal of the Royal Statistical Society Series B-Statistical Methodology*, vol. 67, pp. 301-320, 2005.
- [14] R. Tibshirani, "Regression shrinkage and selection via the Lasso," *Journal of the Royal Statistical Society Series B-Methodological*, vol. 58, pp. 267-288, 1996.
- [15] H. Lei, Z. Huang, J. Zhang, Z. Yang, E.-L. Tan, F. Zhou, *et al.*, "Joint detection and clinical score prediction in Parkinson's disease via multi-modal sparse learning," *Expert Systems with Applications*, vol. 80, pp. 284-296, Sep 2017.
- [16] F. Liu, C.-Y. Wee, H. Chen, and D. Shen, "Inter-modality relationship constrained multi-modality multi-task feature selection for Alzheimer's Disease and mild cognitive impairment identification," *NeuroImage*, vol. 84, pp. 466-475, 2014.
- [17] X. He, D. Cai, and P. Niyogi, "Laplacian score for feature selection," in *the Proceedings of the 18th International Conference on Neural Information Processing Systems*, 2005, pp. 507-514.
- [18] I. Jolliffe, *Principal component analysis*: Springer, 2011.
- [19] L. Shi, L. Du, and Y. D. Shen, "Robust Spectral Learning for Unsupervised Feature Selection," in *2014 IEEE International Conference on Data Mining*, 2014, pp. 977-982.
- [20] S. Wang, J. Tang, and H. Liu, "Embedded Unsupervised Feature Selection," in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015, pp. 470-476.
- [21] X. Zhu, H. I. Suk, and D. Shen, "A novel matrix-similarity based loss function for joint regression and classification in AD diagnosis," *Neuroimage*, vol. 100, pp. 91-105, Oct 2014.
- [22] X. Zhu, H.-I. Suk, L. Wang, S.-W. Lee, and D. Shen, "A novel relational regularization feature selection method for joint regression and classification in AD diagnosis," *Medical Image Analysis*, vol. 38, pp. 205-214, May 2017.
- [23] L. Liu, Q. Wang, E. Adeli, L. Zhang, H. Zhang, and D. Shen, "Exploring diagnosis and imaging biomarkers of Parkinson's disease via iterative canonical correlation analysis based feature selection," *Computerized Medical Imaging and Graphics*, vol. 67, pp. 21-29, 2018.
- [24] E. Adeli, G. Wu, B. Saghafi, L. An, F. Shi, and D. Shen, "Kernel-based Joint Feature Selection and Max-Margin Classification for Early Diagnosis of Parkinson's Disease," *Scientific Reports*, vol. 7, p. 41069, 2017.
- [25] P. A. Pérez-Toro, J. C. Vázquez-Correa, T. Arias-Vergara, E. Nöth, and J. R. Orozco-Arroyave, "Nonlinear dynamics and Poincaré sections to model gait impairments in different stages of Parkinson's disease," *Nonlinear Dynamics*, vol. 100, pp. 3253-3276, 2020.
- [26] J. C. Vázquez-Correa, J. R. Orozco-Arroyave, R. Arora, E. Nöth, N. Dehak, H. Christensen, *et al.*, "Multi-view representation learning via gcca for multimodal analysis of Parkinson's disease," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2017, pp. 2966-2970.
- [27] E. Adeli, F. Shi, L. An, C.-Y. Wee, G. Wu, T. Wang, *et al.*, "Joint feature-sample selection and robust diagnosis of Parkinson's disease from MRI data," *NeuroImage*, vol. 141, pp. 206-219, Nov 2016.
- [28] J. Shi, Z. Xue, Y. Dai, B. Peng, Y. Dong, Q. Zhang, *et al.*, "Cascaded Multi-Column RVFL+ Classifier for Single-Modal Neuroimaging-Based Diagnosis of Parkinson's Disease," *IEEE Transactions on Biomedical Engineering*, vol. 66, pp. 2362-2371, 2019.
- [29] P. Yang, F. Zhou, D. Ni, Y. Xu, S. Chen, T. Wang, *et al.*, "Fused Sparse Network Learning for Longitudinal Analysis of Mild Cognitive Impairment," *IEEE Transactions on Cybernetics*, pp. 1-14, 2019.
- [30] D. Zhang and D. Shen, "Multi-modal multi-task learning for joint prediction of multiple regression and classification variables in Alzheimer's disease," *Neuroimage*, vol. 59, pp. 895-907, Jan 2012.
- [31] B. Lei, P. Yang, T. Wang, S. Chen, and D. Ni, "Relational-Regularized Discriminative Sparse Learning for Alzheimer's Disease Diagnosis," *IEEE Trans Cybern*, vol. 47, pp. 1102-1113, Apr 2017.
- [32] D. Cai, C. Zhang, and X. He, "Unsupervised feature selection for multi-cluster data," in *the Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2010, pp. 333-342.
- [33] F. Nie, D. Xu, I. W. H. Tsang, and C. Zhang, "Flexible Manifold Embedding: A Framework for Semi-Supervised and Unsupervised Dimension Reduction," *IEEE Transactions on Image Processing*, vol. 19, pp. 1921-1932, 2010.
- [34] C. Hou, F. Nie, X. Li, D. Yi, and Y. Wu, "Joint Embedding Learning and Sparse Regression: A Framework for Unsupervised Feature Selection," *IEEE Transactions on Cybernetics*, vol. 44, pp. 793-804, 2014.
- [35] B. Mohar, Y. Alavi, G. Chartrand, O. Oellermann, and A. Schwenk, "The Laplacian spectrum of graphs," *Graph theory, combinatorics, and applications*, vol. 2, pp. 871-898, 1991.
- [36] K. Fan, "On a Theorem of Weyl Concerning Eigenvalues of Linear Transformations. I," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 35, pp. 652-655, 1949.
- [37] Y. C. Pai and J. Patton, "Center of mass velocity-position predictions for balance control," *J Biomech*, vol. 30, pp. 347-54, Apr 1997.
- [38] S. A. Sathanathan, W. Zheng, M. W. Chee, and V. Zagorodnov, "Skull stripping using graph cuts," *Neuroimage*, vol. 49, pp. 225-239, Jan 2010.
- [39] N. Tzourio-Mazoyer, B. Landeau, D. Papathanassiou, F. Crivello, O. Etard, N. Delcroix, *et al.*, "Automated anatomical labeling of activations in SPM using a macroscopic anatomical parcellation of the MNI MRI single-subject brain," *Neuroimage*, vol. 15, pp. 273-289, Jan 2002.
- [40] M. Jenkinson, C. F. Beckmann, T. E. Behrens, M. W. Woolrich, and S. M. Smith, "FSL," *Neuroimage*, vol. 62, pp. 782-790, Aug 2012.
- [41] M. Xia, J. Wang, and Y. He, "BrainNet Viewer: A Network Visualization Tool for Human Brain Connectomics," *PLOS ONE*, vol. 8, p. e68910, 2013.
- [42] R. Gonzalez-Redondo, D. Garcia-Garcia, P. Clavero, C. Gasca-Salas, R. Garcia-Eulate, J. L. Zubieta, *et al.*, "Grey matter hypometabolism and atrophy in Parkinson's disease with cognitive impairment: a two-step process," *Brain*, vol. 137, pp. 2356-67, Aug 2014.
- [43] D. Zhang, J. Wang, X. Liu, J. Chen, and B. Liu, "Aberrant Brain Network Efficiency in Parkinson's Disease Patients with Tremor: A Multi-Modality Study," *Frontiers in aging neuroscience*, vol. 7, pp. 169-169, 2015.
- [44] C. Nombela, L. E. Hughes, A. M. Owen, and J. A. Grahn, "Into the groove: can rhythm influence Parkinson's disease?," *Neurosci Biobehav Rev*, vol. 37, pp. 2564-2570, Dec 2013.

- [45] T. Wu and M. Hallett, "A functional MRI study of automatic movements in patients with Parkinson's disease," *Brain*, vol. 128, pp. 2250-2259, 2005.
- [46] J. Kim, M. Criaud, S. S. Cho, M. Díez-Cirarda, A. Mihaescu, S. Coakeley, *et al.*, "Abnormal intrinsic brain functional network dynamics in Parkinson's disease," *Brain*, vol. 140, pp. 2955-2967, 2017.
- [47] N. Tuovinen, K. Seppi, F. de Pasquale, C. Müller, M. Nocker, M. Schocke, *et al.*, "The reorganization of functional architecture in the early-stages of Parkinson's disease," *Parkinsonism & Related Disorders*, vol. 50, pp. 61-68, 2018.
- [48] R. González-Redondo, D. García-García, P. Clavero, C. Gasca-Salas, R. García-Eulate, J. L. Zubieta, *et al.*, "Grey matter hypometabolism and atrophy in Parkinson's disease with cognitive impairment: a two-step process," *Brain*, vol. 137, pp. 2356-2367, 2014.
- [49] Y. Hosokai, Y. Nishio, K. Hirayama, A. Takeda, T. Ishioka, Y. Sawada, *et al.*, "Distinct patterns of regional cerebral glucose metabolism in Parkinson's disease with and without mild cognitive impairment," *Movement Disorders*, vol. 24, pp. 854-862, 2009.
- [50] J. Xu, F. Jiao, Y. Huang, X. Luo, Q. Xu, L. Li, *et al.*, "A Fully Automatic Framework for Parkinson's Disease Diagnosis by Multi-Modality Images," *Frontiers in Neuroscience*, vol. 13, 2019.
- [51] M. Ariz, R. C. Abad, G. Castellanos, M. Martínez, A. Muñoz-Barrutia, M. A. Fernández-Seara, *et al.*, "Dynamic Atlas-Based Segmentation and Quantification of Neuromelanin-Rich Brainstem Structures in Parkinson Disease," *IEEE Transactions on Medical Imaging*, vol. 38, pp. 813-823, 2019.
- [52] H. Chen, Q. Dou, L. Yu, J. Qin, and P.-A. Heng, "VoxResNet: Deep voxelwise residual networks for brain segmentation from 3D MR images," *NeuroImage*, vol. 170, pp. 446-455, 2018.
- [53] X. Wang, D. Shen, and H. Huang, "Prediction of Memory Impairment with MRI Data: A Longitudinal Study of Alzheimer's Disease," in *Medical Image Computing and Computer-Assisted Intervention*, 2016, pp. 273-281.
- [54] C. Klein and A. Westenberger, "Genetics of Parkinson's disease," *Cold Spring Harbor perspectives in medicine*, vol. 2, p. a008888, 2012.
- [55] D. Zhu, K. Li, L. Guo, X. Jiang, T. Zhang, D. Zhang, *et al.*, "DICCCOL: dense individualized and common connectivity-based cortical landmarks," *Cereb Cortex*, vol. 23, pp. 786-800, Apr 2013.



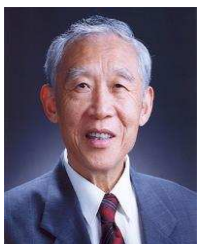
Zhongwei Huang received the M.Eng. degree from Shenzhen University, Shenzhen, China, in 2019. He is also currently pursuing the Ph.D. degree at Shenzhen University.

His current research interests include medical image analysis, machine learning, and pattern recognition.



Haijun Lei received the Ph.D. degree from Huazhong University of Science and Technology, Wuhan, in 2003.

He is currently with the College of Computer Science and Software Engineering, Shenzhen University, China. His current research interests include medical image analysis, machine learning, and pattern recognition.



Guoliang Chen received the B.Sc. degree from Xi'an Jiaotong University, Xi'an, China, in 1961.

His research interests include parallel algorithm, computer architecture, computer network, and computational intelligence.

Prof. Chen is an Academician of Chinese Academy of Sciences.



Alejandro F. Frangi received Ph.D. degree from Utrecht University, Utrecht, in 2001.

His current research interests include medical image computing, medical imaging, and image-based computational physiology.



Yanwu Xu received the B.Eng. and PhD degrees from the University of Science and Technology of China, in 2004 and 2009, respectively.

His current research interests include ocular imaging, medical image analysis, computer vision, and machine learning.



Ahmed Elazab received the Ph.D. degree from University of Chinese Academy of Sciences, China, in 2017.

His current research interests include machine and deep learning, pattern recognition, medical image analysis, and computer-aided diagnosis.



Jing Qin received the Ph.D. degree from The Chinese University of Hong Kong, Hong Kong, in 2009.

His current research interests include innovations for healthcare and medicine applications, medical image processing, deep learning, visualization and human-computer interaction, and health informatics.



Baiying Lei received her M.Eng. degree in electronics science and technology from Zhejiang University, China, in 2007, and Ph.D. degree from Nanyang Technological University (NTU), Singapore, in 2013.

She is currently with School of Biomedical Engineering, Health Science Center, Shenzhen University, China. Her current research interests include medical image analysis, machine learning, and pattern recognition.

Dr. Lei has coauthored more than 120 scientific articles, e.g., IEEE TCYB, IEEE TMI, IEEE TBME, IEEE JBHI. Pattern Recognition and Information Sciences. She is an IEEE senior member and serves as the editorial board member of Scientific Reports, Frontiers in Neuroinformatics, Frontiers in Aging Neuroscience, and Academic Editor of Plos One.