

Early Identification of Mild Cognitive Impairment Using Incomplete Random Forest-Robust Support Vector Machine and FDG-PET Imaging

Shen Lu¹, Yong Xia², Weidong Cai¹, David Dagan Feng¹, and Michael Fulham^{3,4}, for the Alzheimer's Disease Neuroimaging Initiative*

¹ Biomedical and Multimedia Information Technology (BMIT) Research Group, School of Information Technologies, University of Sydney, Sydney, Australia

² Shaanxi Key Lab of Speech & Image Information Processing (SAIIP), School of Computer Science, Northwestern Polytechnical University, Xi'an, China

³ Department of PET and Nuclear Medicine, Royal Prince Alfred Hospital (RPAH), Sydney, Australia

⁴ Sydney Medical School, University of Sydney, Sydney, Australia

*Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Ma-nuscript_Citations.pdf

ABSTRACT

Alzheimer's disease (AD) is the most common type of dementia and will be an increasing health problem in society as the population ages. Mild cognitive impairment (MCI) is considered to be a prodromal stage of AD. The ability to identify subjects with MCI will be increasingly important as disease modifying therapies for AD are developed. We propose a semi-supervised learning method based on robust optimization for the identification of MCI from $[^{18}F]Fluorodeoxyglucose$ PET scans. We extracted three groups of spatial features from the cortical and subcortical regions of each FDG-PET image volume. We measured the statistical uncertainty related to these spatial features via transformation using an incomplete random forest and formulated the MCI identification problem under a robust optimization framework. We compared our approach to other state-of-the-art methods in different learning schemas. Our method outperformed the other techniques in the ability to separate MCI from normal controls.

Keywords: FDG-PET, Alzheimer's disease, Mild Cognitive Impairment, robust optimization

1. INTRODUCTION

Alzheimer's disease (AD) is a neurodegenerative brain disorder that is characterized by progressive memory loss, cognitive impairment and the inability to perform usual daily activities [1]. It is the most common type of dementia, accounting for about 65% of all dementia cases globally and the number of patients is increasing every year as people live longer [2]. Mild cognitive impairment (MCI) is considered as the prodromal phase of AD [3]. Individuals with MCI show greater cognitive impairment than expected for their age, but they do not meet the criteria for dementia [4]. The conversion rate of MCI to AD is estimated to be between 10% - 25% per year [5]. Although there are no current disease modifying agents to halt the progression of AD there are a number of clinical trials underway in patients with pre-symptomatic disease [6]. Thus as effective therapies become available the early identification of patients with MCI will be of tremendous benefit to patients and their families.

The pathology of AD includes cortical and subcortical atrophy together with the deposition of β -amyloid. Two widely used AD biomarkers are structural imaging with magnetic resonance (MR) imaging [7] and functional imaging with [^{18}F]Fluorodeoxyglucose positron emission tomography (FDG-PET) [2]. The advantage of FDG-PET over MR imaging is that PET can detect reduced cerebral glucose metabolism before structural change is evident on MR imaging. The separation of patients with MCI from normal controls (NCs) by the visual analysis of FDG-PET images, however, is difficult. Visual interpretation of these studies is also operator-dependent and related to the skill and experience of the reader. A reliable and robust computer-aided method could improve this situation.

Machine learning theory has been applied to the dementias and Davatzikos et al. used a voxel-based nonlinear multivariate analysis to separate AD from Frontotemporal dementia (FTD) with MR imaging [8]. In their subsequent study [9], they applied a similar method to combinations of features extracted from MR images and cerebrospinal fluid (CSF) to predict progression from MCI to AD. Although there are a number of pathological studies of MCI with PET [2, 10], the use of computerized classification methods based on PET data is not prominent in the literature. In a previous study [11], we implemented a method that combined multi-kernel learning (MKL) and a genetic algorithm (GA) to differentiate between AD, FTD and NC with FDG-PET images. We used GA to obtain the optimal kernel weights for combining different kernel matrices and then trained a MKL machine to classify the three classes at the same time. In a subsequent study [12], we used an automated classification method for dealing with AD and NC using infinite kernel learning (IKL). We exploited the importance of cerebral features in the AD/NC

classification task using this method. We investigated the early identification of different dementia sub-types using FDG-PET and reported superior classification accuracy and efficiency, but we did not address the problem of separating MCI and NCs. Zhang et al. in [13] combined a number of biomarkers (MR, PET and CSF) together and used MKL to classify AD, MCI and NC. They reported good differentiation of AD from NC but they had a lower accuracy (76.4%) for separating MCI from NCs. In addition, in the clinical setting it is difficult to obtain all three biomarkers due to costs and the reluctance for subjects to undertake a lumbar puncture. Recently, Gray et al. [14] proposed a multi-modality classification process based on the embedding of feature similarities among MR, FDG-PET, CSF, and genetic information via random forest (RF). They reported 75% classification accuracy between MCI and NCs which was poorer than the 89% accuracy in separating AD from NCs.

In this work our aim was the early identification of patients with MCI using FDG-PET imaging. We used an incomplete random forest - robust support vector machine (IRF-RSVM) approach to address the problem where subjects with MCI have similar imaging to NCs and the spatial resolution of FDG-PET is poorer than structural imaging. The idea was to build an incomplete random forest using FDG-PET image features and model the outputs of the random forest as a noise corrupted feature dataset, and then minimize a loss function in terms of these noisy data within a robust programming framework.

2. BACKGROUND

2.1 Random Forest (RF)

Random forest is an ensemble learning method, which builds a number of decision trees [15, 16] with random factors. Basically, RF injects randomness into its learning process in two forms: random sampling and random parameterization. Random sampling arbitrarily selects training examples to train each decision tree. Random parameterization chooses training parameters during the training of each decision tree in an unplanned fashion. Both or either of these two forms of randomness can be used in the training process. The introduced randomness prompts variation and diversity among the decision trees that are built. Each decision tree in the forest is a binary tree on which each non-leaf node, a so-called weak learner, is trained by solving an optimization problem to determine the best data feature to use to split the dataset. For features with a numerical value, we simply threshold the data set at the current node so that examples, where the value of the feature used for splitting is less than the threshold, go to the

left branch of this node and other examples go to the right branch. The process continues on subsequent nodes until a stopping criterion is met.

2.2 Support Vector Machine

In general, the goal of machine learning is to learn distinguishable patterns from training data belonging to different classes, and then use these patterns to classify new (unseen) data (test data) to some extent. Kernel based maximum margin learning methods have been very widely used in machine learning research during the last decade [17-19]. Basically, kernel based method constructs kernels in reproducing kernel Hilbert space (RKHS) based on data, and finds a separating hyperplane that separates data belonging to different classes with maximum margins by minimizing a structural empirical risk functional [17, 19]. Within this family of methods, support vector machine (SVM) is the most well-known method and has been used in many scientific and industrial applications [17].

SVM finds the optimal separating hyperplane by solving a linearly constrained quadratic optimization problem (QP), which can be written as:

$$\text{minimize}_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{k=1}^p \xi_k \quad (1)$$

$$\text{s. t. } y_k(\mathbf{w}^T \mathbf{x}_k + b) \geq 1 - \xi_k, \quad \xi_k \geq 0, \quad \forall k = 1, \dots, p$$

where \mathbf{x} is the training data vector with label $y \in \{-1, 1\}$, p is the number of training data, ξ is slack variable which allows some data to be misclassified, the weight vector \mathbf{w} and bias b are optimization variables that define the hyperplane. Solving the optimization problem (1) results in a separating hyperplane that separates training data, and at the same time, maximizes the margins between training data on both sides of the hyperplane [17]. After solving (1), the prediction of testing data label is made by evaluating the function below for each testing data \mathbf{x}' :

$$f(\mathbf{x}'; \mathbf{w}^*, b^*) = \text{sgn}(\mathbf{w}^{*T} \mathbf{x}' + b^*) \quad (2)$$

where \mathbf{w}^* and b^* are the optimal solutions of (1), $\text{sgn}(\cdot)$ gives the sign of the operator and the sign indicates the class membership of testing data \mathbf{x}' .

2.3 Inductive Learning and Transductive Learning

Theoretically, there are two types of machine learning schemas, inductive and transductive learning. In the

inductive learning setting, a learner is trained using a set of observed data called training data and is then tested on a set of previously unseen data called test data. This setting is extremely common in machine learning research. Transductive learning differs from inductive learning in that, during the training phase a participant has visibility of training data and test data and a participant can potentially make use of the information, exposed by the test data, such as the probability distribution information [19, 20]. Hence, transductive learning is ideal when the size of the experimental data is small. In this work we tested the proposed method in the inductive and transductive learning settings.

3. DATA AND MATERIALS

The FDG-PET image data we used were from the Alzheimer’s Disease Neurodegenerative Initiative (ADNI) cohort (<http://adni.loni.usc.edu>). ADNI is a multi-center program funded by a public-private partnership and non-profit organizations to provide standardized longitudinal medical image data to global researchers for neurodegenerative disease research. There were 140 FDG-PET studies; 70 MCI subjects and 70 NCs. All images came from ADNI, ADNI GO and ADNI 2 baseline/initial scans; these data had been through a pre-processing pipeline that included: co-registration, averaging, voxel normalization, and isotropic Gaussian smoothing [21]. This pre-processing work is done by the ADNI participants and it makes any subsequent analysis simpler as the data from different PET scanners are then uniform. The demographic information of all 140 subjects and the Mini-Mental State Examination (MMSE) scores are shown in Table I.

TABLE I

Demographics data with mean \pm std

	Subjects	
	MCI	NC
Number of subjects	70	70
Age (years)	75.5 \pm 7.4	74.5 \pm 5.7

Weight (kg)	77.9±14.8	75.5±15.2
Gender (M/F)	45/25	42/28
MMSE	26.8±1.6	28.92±1.3

4. METHODS

4.1 Feature Extraction

Our aim was to extract spatial features from voxel volumes representing cerebral cortical and subcortical regions on each PET image. To ensure good spatial localization we registered each PET image to a brain atlas. We used the automated anatomical labeling (AAL) cortical parcellation map [22] to identify the anatomical volumes of interest (VOIs) where spatial features were to be extracted. Hanning et al. [10] reported on the important role that the AAL map plays in computer-based functional brain image analysis for identifying dementia. The AAL image template contains 116 manually drawn and accurately reconstructed anatomical VOIs, and it has dimension of $91 \times 109 \times 91$ with voxel size of $2 \times 2 \times 2 \text{ mm}^3$. To achieve the best image registration result, before registration to AAL template we spatially normalized each of the study images to the PET image template provided by statistical parametric mapping (SPM) software. This PET template has the same dimension and voxel size as the AAL template. As a result, the normalized images are in the same coordinate space as the AAL template.

Following spatial normalization, we extracted three groups of spatial features. They were: mean voxel values from 116 anatomical VOIs, standard deviations of voxel values from 116 anatomical VOIs, and mean voxel value differences between 54 pairs of anatomical VOIs on left and right brain hemispheres. We then concatenated these three feature groups together to form a feature vector of dimension 286 for each image.

Let $\mathbf{X} \in \mathbb{R}^{140 \times 286}$ denote the column matrix containing all spatial features, and let $\mathbf{x}_i = \{\mathbf{x}_{i,1}, \mathbf{x}_{i,2}, \dots, \mathbf{x}_{i,286}\}^T$, $i = 1, \dots, 140$ be the feature vector for the i th image, where T is matrix transpose. Finally, let $\mathbf{Y} = \{-1, 1\}$, $\mathbf{Y} \in \mathbb{R}^{140}$ denote the label vector (MCI: -1 , NC: 1) for the images. Note that \mathbf{X} is mean centered and standardized.

4.2 Feature Transformation

In our method we do not use the feature matrix \mathbf{X} directly. Instead, we used a transformed version of \mathbf{X} , because

>CMIG<

we wanted to better model the classification problem with noise corrupted images in the robust optimization framework. We attempted to model image noise caused by perturbation to the data \mathbf{X} as a perturbation to the statistical distribution of \mathbf{X} . Therefore, we use $\tilde{\mathbf{X}}$ to denote the transformed data matrix, and $\tilde{\mathbf{Y}}$ the transformed label vector.

The data transformation process took the form of incomplete random forest, whose main difference compared to the classic random forest is that the decision trees in the incomplete random forest are never fully grown. That is to say, the training cycle of each decision tree is terminated before it reaches the state when each leaf tree node contains only data examples from single class. Let \mathcal{F} be the incomplete random forest we build and denote the decision tree as T_m , $m = 1, 2, \dots, N_T$ where N_T is the total number of trees in \mathcal{F} . T_m is only allowed to grow up to d level where d is a predefined parameter.

To construct \mathcal{F} , we iteratively built each decision tree T_m using \mathbf{X} by branching \mathbf{X} at each non-leaf tree node following top-down order. Starting from the root node of T_m , at each non-leaf node of the tree we randomly selected a number of different features $k = 1, \dots, n_k$ (where n_k is a predefined parameter) from the 286 spatial features and calculate the branching threshold θ_k using

$$\theta_k = (\max(\mathbf{x}_{:,k}) - \min(\mathbf{x}_{:,k})) / 2 \quad (3)$$

where \max and \min indicate the maximum and minimum values in a vector, respectively. $\mathbf{x}_{:,k} \in \mathbb{R}^{140}$ is the column vector of k th feature values in $\mathbf{x}_i, i = 1, \dots, 140$. We then selected the best branching threshold θ^* from the candidates $\theta_1, \dots, \theta_{n_k}$ by solving the optimization problem below

$$\theta^* = \operatorname{argmin}_{\theta_k} I \quad k = 1, \dots, n_k \quad (4)$$

where I is the unsupervised information gain [15] defined by

$$I = \log(|\Lambda(S_h)|) - \sum_{b=\{L,R\}} |S_h^b| \log(|\Lambda(S_h^b)|) / |S_h| \quad (5)$$

where h is the current non-leaf node being branched, $S_h \subset \mathbf{X}$ is the dataset at node h before branching, b is the branching direction which can only be either L - left branch of node h or R - right branch of node h , $S_h^b \subset S_h$ is thus the dataset assigned to the respective branch (left or right) of node h , Λ is covariance operator and $|\cdot|$ is set cardinality. Note that the calculated unsupervised information gain I may not be a real number due to the presence of

>CMIG<

the covariance operator, in which case the candidate θ_k is discarded and a new θ is randomly selected to replace it. This branching/optimization process is carried on until the predefined tree depth d is reached.

Once every T_m in \mathcal{F} is built following the procedure outlined above, the transformed feature data and labels are collected from leaf nodes of each tree T_m . Each leaf node of T_m is treated as a subspace containing two clusters, one for each of the two data classes (MCI, NC). It is straightforward to calculate the mean μ and covariance Σ from each cluster. For T_m we obtain $\boldsymbol{\mu}_m = \{\mu_+^1, \dots, \mu_+^l, \mu_-^1, \dots, \mu_-^l\}^T$, $\boldsymbol{\Sigma}_m = \{\Sigma_+^1, \dots, \Sigma_+^l, \Sigma_-^1, \dots, \Sigma_-^l\}^T$, and $\mathbf{y}_m = \{1, \dots, 1, -1, \dots, -1\}^T$, $|\mathbf{y}_m| = |\boldsymbol{\mu}_m| = |\boldsymbol{\Sigma}_m|$ where l is the number of leaf nodes, $+$ is MCI class, $-$ is NC class, \mathbf{y} is the corresponding label vector. To this end, the transformed data is

$$\tilde{\mathbf{X}} = \{\mathcal{S}_\mu, \mathcal{S}_\Sigma\} = \{\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m\}, \{\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_m\}\}$$

and the transformed label vector is

$$\tilde{\mathbf{Y}} = \{\mathbf{y}_1, \dots, \mathbf{y}_m\}$$

4.3 Classification

In the classification stage, we train a classifier within the robust optimization (RO) framework [23] using transformed data $\{\mathcal{S}_\mu, \mathcal{S}_\Sigma, \tilde{\mathbf{Y}}\}$. Assuming that $\{\mathcal{S}_\mu, \mathcal{S}_\Sigma, \tilde{\mathbf{Y}}\}$ are noise corrupted, which is appropriate as it is accepted that FDG-PET images usually have a low signal-to-noise ratio due to the limited resolution of PET scanners [24], to model the uncertainty associated with these noisy data, we consider the modified version of the inequality constraint in the original SVM problem

$$Pr\{y_k(\mathbf{w}^T \mathbf{x}_k + b) \geq 1 - \xi_k\} \geq \delta_k \quad (6)$$

where $\delta \in [0,1)$ is a user defined parameter. The probabilistic constraint simply requires each feature vector \mathbf{x}_k to lie on the correct side of the optimal hyperplane with a certain confidence value δ . Solving SVM problem with this constraint is extremely difficult. Therefore, we transformed it into a deterministic constraint with the assumption that the feature data is drawn from a multimodal Gaussian distribution characterized by mean and covariance [25]. Our transformed datasets $\{\mathcal{S}_\mu, \mathcal{S}_\Sigma, \tilde{\mathbf{Y}}\}$ naturally fit into this new deterministic constraint, which is written as

$$\tilde{\mathbf{Y}}(\mathcal{S}_\mu \mathbf{w} + b) \geq 1 - \xi + \gamma \|\mathcal{S}_\Sigma^{1/2} \mathbf{w}\|_2 \quad (7)$$

>CMIG<

where we introduce a new parameter vector $\boldsymbol{\gamma}$, $|\boldsymbol{\gamma}| = |\tilde{\mathbf{Y}}|$. $\boldsymbol{\gamma}$ is computed from the leaf nodes of decision trees in a way similar to [26]. For $i = 1, \dots, lm$, where l is the number of the number of leaf nodes, m is the number of trees in forest \mathcal{F}

$$\gamma_i = \begin{cases} n_i/n_{MCI}, & y_i = 1 \\ n_i/n_{NC}, & y_i = -1 \end{cases} \quad (8)$$

where n_i is the number of feature vectors dwelled at leaf node i , n_{MCI} and n_{NC} are the total number of MCI cases and NC cases in the whole dataset, respectively. Finally, the robust SVM problem is formulated as

$$\text{minimize}_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C\xi \quad (9)$$

$$s. t. \tilde{\mathbf{Y}}(\mathbf{S}_\mu \mathbf{w} + b) \geq 1 - \xi + \boldsymbol{\gamma} \|\mathbf{S}_\Sigma^{1/2} \mathbf{w}\|_2, \quad \xi \geq 0$$

Evidently, (9) is a convex problem. In our study, we solve this problem using CVX Matlab toolbox [27]. In order to efficiently solve this problem with CVX, we reformulate (9) into an equivalent second-order cone programming (SOCP) problem

$$\text{minimize}_{\mathbf{w}, b, \xi, t} t + C\xi \quad (10)$$

$$s. t. \|\mathbf{w}\|_2 \leq t$$

$$\tilde{\mathbf{Y}}(\mathbf{S}_\mu \mathbf{w} + b) \geq 1 - \xi + \boldsymbol{\gamma} \|\mathbf{S}_\Sigma^{1/2} \mathbf{w}\|_2, \quad \xi \geq 0$$

Once the optimal solution $\{\mathbf{w}^*, b\}$ is found by solving (10), predictions of feature vectors extracted from new PET image are made by evaluating function (2).

5. EXPERIMENTS

5.1 Benchmark Methods

We compare the proposed RF-RSVM method to three baseline methods:

1. Supervised SVM [17]: We applied the soft margin SVM as described in the section on background.
2. Laplacian SVM (LapSVM) [28, 29]: LapSVM regularizes the standard SVM cost function with a data dependent penalty term with the assumption that the intrinsic structure of the data is embedded within a low dimensional manifold. It approximates this new penalty term by modeling the structure of the data using graph

Laplacian.

3. Method proposed by Huang et al. [26]: Huang et al conducted clustering based on a dataset using the k-nearest neighbour algorithm, and then merged similar clusters, followed by solving the SOCP problem (10). The method showed good performance on non-medical imaging datasets.

Only the soft margin SVM is supervised learning method while the other two methods are both semi-supervised.

5.2 Experimental Settings

To ensure that our method had good generalizability, we applied 3-fold cross validation for our method and the three benchmark methods. We first divided the whole dataset evenly into 3 subsets (the residual is randomly assigned to one of the subsets), each contained 20% labeled data examples and the rest of data were treated as unlabeled. Within the labeled and unlabeled groups of data in each subset, we further restricted that 50% of data in this group were MCI subjects and 50% were NCs. We used the inductive learning schema in our first experiment, and trained the target classifier using any 2 out of the 3 subsets, then tested the target classifier on the leftover subset. Initially, $d = 3$, $N_T = 25$ were used to construct the unsupervised random forest in RF-RSVM. Hyperparameters required in the benchmark methods were set empirically or selected by an inner 3-fold cross-validation using the training data.

5.3 Results and Discussion

We applied the proposed method (RF-RSVM) and the baseline methods to classify MCI and NC. The performance of these methods measured by classification accuracy, sensitivity and specificity averaged over the 3 fold cross-validation steps are shown in Table II. The plain supervised SVM failed to generate meaningful results possibly due to the high non-linearity and high similarity of the feature patterns in our dataset. The RF-RSVM outperformed the other two semi-supervised learning methods in terms of accuracy, sensitivity and specificity. The less impressive performance of LapSVM may be related to that the intrinsic structure of our data, which to some extent, violates the smooth manifold assumption that is crucial for LapSVM to perform well. The receiver operating characteristic (ROC) curve for these three semi-supervised learning methods are shown in Figure 1 and complement the findings in Table II.

TABLE II

Classification performance of the proposed method and the baseline methods

Method	Accuracy	Sensitivity	Specificity
RF-RSVM	92.18%	90.31%	94.30%
LapSVM	75.54%	83.32%	69.70%
kNN-SVM	86.41%	89.53%	84.20%
Supervised SVM	47.83%	0%	47.83%

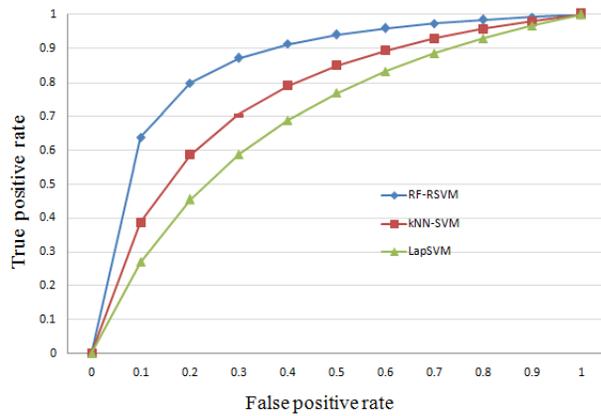


Figure 1. ROC curves for methods, excluding supervised SVM compared in Table I.

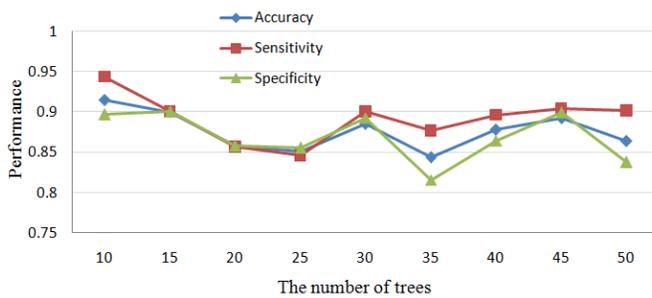
The inductive learning scheme used in the first experiment validates the generalizability of classifier built using training data. Transductive learning, on the other hand, carries out training and testing on the same dataset. It is very useful when the size of the dataset (training + testing) is small. Since the dataset used in dementia related studies usually does not contain tens of thousands images, the proposed method could be tested under the transductive learning setting. For each of the 3 subsets created before, we trained our method, LapSVM, and kNN-SVM in a way that both labeled and unlabeled data within each subset were used for training while only unlabeled data within each subset was used for testing. The same performance metrics as those shown in Table II were used and they were averaged over the 3-fold cross validation. Table III shows the final performance of the three methods.

TABLE III

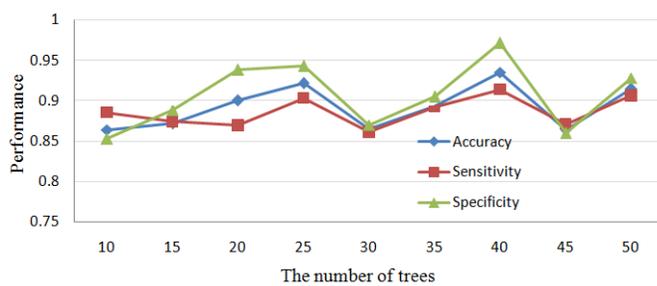
Classification performance of the proposed method and the baseline methods under transductive learning setting

Method	Accuracy	Sensitivity	Specificity
RF-RSVM	86.90%	92.86%	80.95%
LapSVM	77.50%	77.58%	76.67%
kNN-SVM	79.49%	82.35%	77.27%

When dealing with random forest an obvious question is what is the impact that the hyperparameters such as tree depth d and the number trees N_T have on the performance of the proposed method. We carried out experiments to apply RF-RSVM with varying d and N_T to measure the impact. We fixed d as 2 and 3, and then increased the number of trees in the random forest from 10 to 50 with increments of 5. The metrics used were identical to the first experiment and the performance charts for these two scenarios are shown in Figure 2.



(2a)



(2b)

Figure 2. Accuracy, sensitivity, and specificity changes as the number of trees increases with the tree depth fixed as 2 (2a) and 3 (2b).

It is noticeable that our method is not sensitive to the number of trees N_T as the slopes of curves on both charts are relatively stable: on the chart (2a) the performance fluctuates between 0.8 and 0.9, and between 0.85 and 0.95 on the chart (2b). However, increasing d improved the overall performance of the proposed method; and this has been reported for other experiments by Criminisi et al. [15] and Verikas et al. [30]. Criminisi et al. regarded the maximum

>CMIG<

allowed tree depth as one of the most influential factor for a random forest, and Verikas et al. let trees in their random forest grow to maximum depth in order to get low bias and low correlation which are essential for accuracy.

In machine learning, when the whole training dataset is labeled (e.g. each PET image contained in the dataset is given a class: MCI or NC), the learning process is called supervised learning, whereas it is called semi-supervised learning if large part of the training dataset are unlabeled. So if we denote the total number of data examples contained in a dataset as N . Let N_L and N_U be the number of labeled data and unlabeled data. $N = N_L + N_U$, and usually, $N_L \ll N_U$. Therefore, semi-supervised learning can play an important role in solving practical problems when most of data labels are unavailable due to the high cost of manual data labeling or when full data labeling is not possible. Our method is essentially a semi-supervised learning method, which is appropriate, since in the clinical setting brain images labeled as dementia are usually not available given the difficulty in making an accurate diagnosis without a post-mortem. The number of unlabeled brain images, or brain images which are suspected to reflect dementia, are abundant.

One of the most important components/processes in our method is feature transformation via incomplete random forest. This transformation is the key to modeling the MCI/NC classification under RO framework. This transformation also introduces some problems. For example, after a decision tree is constructed it is not guaranteed that each leaf node will always contain some feature vectors belonging to MCI and some belonging to NC. Some leaf nodes may only contain feature vectors belonging to a single class – we call those leaf nodes degenerative leaf nodes. We discard all degenerative leaf nodes to avoid numerical difficulties. The main problem with the feature transformation process is a long training time. This issue can be seen from the first constraint in the optimization problem (10). Recall that d is the depth of each decision tree in forest \mathcal{F} . Since the decision tree is a binary tree, the number of leaf nodes each decision tree can have is $2^d - 2^{d-1}$, thus the total number of leaf nodes in forest \mathcal{F} is $N_T(2^d - 2^{d-1})$. Each leaf node contains two clusters (one for MCI, one for NC), as a result, the upper bound of the number of constraints is $2N_T(2^d - 2^{d-1})$ (this is an upper bound since some leaf nodes may be degenerative). It is easy to have tens of thousands of constraints even with a moderate number of decision trees and tree depth. This greatly decreases the efficiency of our method. A simple strategy to alleviate this effect would be to combine μ and Σ for each data category (MCI and NC) within each tree by calculating their arithmetic means and we applied this strategy to all our experiments.

6. CONCLUSION

We implemented a novel computer-aided method for the early identification of baseline MCI subjects among NCs using FDG-PET image data obtained from the ADNI cohort. We formulated the problem within a robust optimization framework with feature data transformed via incomplete random forest to enable semi-supervised learning. Our results show that our method outperforms two other semi-supervised learning methods. In future work we will test the performance of our method on a much larger dataset to determine if the current results are sustained over a larger dataset.

Conflict of interest statement

The authors do not hold any conflicts of interest that could inappropriately influence this manuscript.

Acknowledgements

This work was supported in part by the National Natural Science Foundation of China under Grants 61471297, in part by the Natural Science Foundation of Shaanxi Province, China, under Grant 2015JM6287, in part by the Fundamental Research Funds for the Central Universities under Grants 3102014JSJ0006, and in part by the Returned Overseas Scholar Project of Shaanxi Province, China.

Data collection and sharing for this project was funded by the Alzheimer's Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and DOD ADNI (Department of Defense award number W81XWH-12-2-0012). ADNI is funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: AbbVie, Alzheimer's Association; Alzheimer's Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-Myers Squibb Company; CereSpir, Inc.; Cogstate; Eisai Inc.; Elan Pharmaceuticals, Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.; Lumosity; Lundbeck; Merck & Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and Transition Therapeutics. The Canadian Institutes of Health Research is providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health (www.fnih.org). The grantee organization is the Northern

California Institute for Research and Education, and the study is coordinated by the Alzheimer's Therapeutic Research Institute at the University of Southern California. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of Southern California.

REFERENCES

- [1] Teune LK, Bartels AL, de Jong BM, Willemsen AT, Eshuis SA, de Vries JJ, et al. Typical cerebral metabolic patterns in neurodegenerative brain diseases. *Mov Disord.* 2010;25:2395-404.
- [2] Devous MD, Sr. Functional brain imaging in the dementias: role in early detection, differential diagnosis, and longitudinal studies. *Eur J Nucl Med Mol Imaging.* 2002;29:1685-96.
- [3] Albert MS, DeKosky ST, Dickson D, Dubois B, Feldman HH, Fox NC, et al. The diagnosis of mild cognitive impairment due to Alzheimer's disease: Recommendations from the National Institute on Aging-Alzheimer's Association workgroups on diagnostic guidelines for Alzheimer's disease. *Alzheimer's & dementia : the journal of the Alzheimer's Association.* 2011;7:270-9.
- [4] McKhann GM, Knopman DS, Chertkow H, Hyman BT, Jack CR, Kawas CH, et al. The diagnosis of dementia due to Alzheimer's disease: Recommendations from the National Institute on Aging-Alzheimer's Association workgroups on diagnostic guidelines for Alzheimer's disease. *Alzheimer's & Dementia.* 2011;7:263-9.
- [5] Grand J, Caspar S, Macdonald S. Clinical features and multidisciplinary approaches to dementia care. *J Multidiscip Healthc.* 2011;4:125-47.
- [6] Morris JC, Aisen PS, Bateman RJ, Benzinger TLS, Cairns NJ, Fagan AM, et al. Developing an international network for Alzheimer's research: the Dominantly Inherited Alzheimer Network. *Clinical Investigation.* 2012;2:975-84.
- [7] Fox NC, Schott JM. Imaging cerebral atrophy: normal ageing to Alzheimer's disease. *The Lancet.* 2004;363:392-4.
- [8] Davatzikos C, Resnick SM, Wu X, Parnpi P, Clark CM. Individual patient diagnosis of AD and FTD via high-dimensional pattern classification of MRI. *NeuroImage.* 2008;41:1220-7.
- [9] Davatzikos C, Bhatt P, Shaw LM, Batmanghelich KN, Trojanowski JQ. Prediction of MCI to AD conversion, via MRI, CSF biomarkers, and pattern classification. *Neurobiology of Aging.* 2011;32:2322.e19-.e27.
- [10] Sun H, Hu B, Yao Z, Jackson M. A PET study of discrimination of cerebral glucose metabolism in Alzheimer's disease and mild cognitive impairment. *Biomedical Engineering and Informatics (BMEI), 2013 6th International Conference on: IEEE; 2013.* p. 6-11.
- [11] Xia Y, Lu S, Wen L, Eberl S, Fulham M, Feng DD. Automated identification of dementia using FDG-PET imaging. *BioMed research international.* 2014;2014.
- [12] Xia Y, Lu S, Wei W, Feng DD, Zhang Y. Non-sparse infinite-kernel learning for automated identification of Alzheimer's disease using PET imaging. *Control Automation Robotics & Vision (ICARCV), 2014 13th International Conference on: IEEE; 2014.* p. 855-60.
- [13] Zhang D, Wang Y, Zhou L, Yuan H, Shen D. Multimodal classification of Alzheimer's disease and mild cognitive impairment. *NeuroImage.* 2011;55:856-67.
- [14] Gray KR, Aljabar P, Heckemann RA, Hammers A, Rueckert D. Random forest-based similarity measures for multi-modal classification of Alzheimer's disease. *NeuroImage.* 2013;65:167-75.
- [15] Criminisi A, Shotton J, Konukoglu E. Decision forests: A unified framework for classification, regression, density estimation, manifold learning and semi-supervised learning. *Foundations and Trends® in Computer Graphics and Vision.* 2012;7:81-227.
- [16] Breiman L. Random forests. *Machine learning.* 2001;45:5-32.
- [17] Schölkopf B, Smola AJ. *Learning with kernels : support vector machines, regularization, optimization, and beyond.* Cambridge, Mass.: MIT Press; 2002.
- [18] Xu L, Neufeld J, Larson B, Schuurmans D. Maximum margin clustering. *Advances in neural information processing systems2004.* p. 1537-44.
- [19] Vapnik V. *The nature of statistical learning theory:* Springer Science & Business Media; 2013.
- [20] Joachims T. Transductive learning via spectral graph partitioning. *ICML2003.* p. 290-7.
- [21] ADNI. PET Pre-processing.
- [22] Tzourio-Mazoyer N, Landeau B, Papathanassiou D, Crivello F, Etard O, Delcroix N, et al. Automated anatomical labeling of activations in SPM using a macroscopic anatomical parcellation of the MNI MRI single-subject brain. *NeuroImage.* 2002;15:273-89.
- [23] Ben-Tal A, Nemirovski A. Robust convex optimization. *Mathematics of operations research.* 1998;23:769-805.
- [24] Feng DD. *Biomedical information technology:* Academic Press; 2011.
- [25] Shivaswamy PK, Bhattacharyya C, Smola AJ. Second order cone programming approaches for handling missing and uncertain data. *The Journal of Machine Learning Research.* 2006;7:1283-314.

- [26] Huang G, Song S, Gupta JN, Wu C. A second order cone programming approach for semi-supervised learning. *Pattern Recognition*. 2013;46:3548-58.
- [27] Grant M, Boyd S, Ye Y. CVX: Matlab software for disciplined convex programming. 2008.
- [28] Belkin M, Niyogi P, Sindhvani V. Manifold Regularization: A Geometric Framework for Learning from Labeled and Unlabeled Examples. *J Mach Learn Res*. 2006;7:2399-434.
- [29] Melacci S, Belkin M. Laplacian support vector machines trained in the primal. *The Journal of Machine Learning Research*. 2011;12:1149-84.
- [30] Verikas A, Gelzinis A, Bacauskiene M. Mining data with random forests: A survey and results of new tests. *Pattern Recognition*. 2011;44:330-49.

LIST OF TABLES AND FIGURES

TABLE I

Summary Of Demographics of Data set with mean \pm std

TABLE II

Classification performance of the proposed method and the baseline methods

TABLE III

Classification performance of the proposed method and the baseline methods under transductive learning setting

Figure I. ROC curves for methods (excl. supervised SVM) compared in TABLE I.

Figure II. Accuracy, sensitivity, and specificity changes as the number of trees increases with the tree depth fixed as 2 (2a) and 3 (2b).