

Diagnosis of Schizophrenia: A comprehensive evaluation

M. Tanveer^a, Jatin Jangir^b, M.A. Ganaie^a, Iman Beheshti^c, M. Tabish^a, Nikunj Chhabra^b

^a*Department of Mathematics, Indian Institute of Technology Indore, Simrol, Indore, 453552, India*

^b*Department of Electrical Engineering, Indian Institute of Technology Indore, Simrol, Indore, 453552, India*

^c*Department of Human Anatomy and Cell Science, Rady Faculty of Health Sciences, Max Rady College of Medicine, University of Manitoba, Winnipeg, MB, Canada.*

Abstract

Machine learning models have been successfully employed in the diagnosis of Schizophrenia disease. The impact of classification models and the feature selection techniques on the diagnosis of Schizophrenia have not been evaluated. Here, we sought to access the performance of classification models along with different feature selection approaches on the structural magnetic resonance imaging data. The data consist of 72 subjects with Schizophrenia and 74 healthy control subjects. We evaluated different classification algorithms based on support vector machine (SVM), random forest, kernel ridge regression and randomized neural networks. Moreover, we evaluated T-Test, Receiver Operator Characteristics (ROC), Wilcoxon, entropy, Bhattacharyya, Minimum Redundancy Maximum Relevance (MRMR) and Neighbourhood Component Analysis (NCA) as the feature selection techniques. Based on the evaluation, SVM based models with Gaussian kernel proved better compared to other classification models and Wilcoxon feature selection emerged as the best feature selection approach. Moreover, in terms of data modality the performance on integration of the grey matter and white matter proved better compared to the performance on the grey and white matter individually. Our evaluation showed that classification algorithms along with the feature selection approaches impact the diagnosis of Schizophrenia disease. This indicates that proper selection of the features and the classification models can improve the diagnosis of Schizophrenia.

Keywords: Schizophrenia, classification, machine learning.

1. Introduction

Schizophrenia is a severe mental disorder that affects millions of people worldwide. Schizophrenia makes people slowly lose contact with reality, leading to hallucinations, delusions, and extremely disordered thinking. Patients report hearing voices or seeing things that are not there; they also tend to develop fixed and false beliefs. Suicidal tendency is also

Email addresses: mtanveer@iiti.ac.in (M. Tanveer), ee180002022@iiti.ac.in (Jatin Jangir), phd1901141006@iiti.ac.in (M.A. Ganaie), Iman.beheshti@umanitoba.ca (Iman Beheshti), mcs1903141002@iiti.ac.in (M. Tabish), ee180002040@iiti.ac.in (Nikunj Chhabra)

Preprint submitted to IEEE

March 23, 2022

a common trait among Schizophrenia patients. Moreover, patients inflicted with Schizophrenia are 2 – 3 times more likely to die than the general public due to patients not seeking aid for preventable physical diseases [16]. Luckily, Schizophrenia is treatable with medicines and psycho-social support, and these methods have proven successful [37]. Thus, the central blockade in eradicating Schizophrenia is lack of its early detection.

Several attempts have been made to remedy this problem. Many studies show promising results but, Machine Learning (ML) has seen little use in clinical practice for Schizophrenia. This can perhaps be credited to the unreliability or stability of some machine learning models; this influences the consensus that ML solutions are not dependable and cannot be trusted, especially for a medical job. Even though doctors sometimes make mistakes themselves, people still trust them. Developing this same level of trust for a machine would be an arduous task. Machine learning should be seen as a research tool to advance the field of study, not the be-all and end-all. In the past few years, extensive research has been done on various classification algorithms and their improved versions have been proposed like for SVM, it's extensions like twin support vector machine (TWSVM) [9], twin bounded SVM (TBSVM) [25] etc. are proposed to improve the performance of SVM algorithm. Other methods such as k -nearest neighbour (KNN) [7], random forest (RaF) [1] have also been thoroughly studied. Interested readers can refer to the comprehensive review on TWSVM [32].

Our intent with this study is to perform a comprehensive evaluation and bring awareness towards using different classification algorithms and their variants and extensions for diagnosing schizophrenia disease. In this study, we evaluate single-modal methods using exclusively Structural MRI (sMRI) scans to train and validate them. We use the same dataset for all the algorithms. The results of this study will help choose a suitable classification algorithm and feature selection technique based on the requirement. The rest of the paper is organized as follows: In Section 3, we discuss about subjects, 3D MRI processing and give a brief description of various classification algorithms and feature selection techniques used and also about validation and experimental setup. Performance of various classification algorithms on white matter, grey matter and integrated matter is discussed in Section 4. Analysis and summarizing of results is done in Section 5. Conclusions and future works are discussed in Section 6. Please note, Figures and Tables referenced from the attached Supplementary Paper have been suffixed with an "S".

2. Related Works

The current diagnosis scheme for Schizophrenia is to rule out other mental disorders and then employ psychiatric and physical screening. The studies [26, 27] evaluated magnetic resonance imaging (MRI) scans for the detection of Schizophrenia. A review of MRI findings in schizophrenia [26] discusses brain abnormalities due to Schizophrenia. Thus, there is significant evidence that by using MRI scans, one can exploit ML techniques to automate and improve the detection of Schizophrenia. Several studies [5, 39] have already attempted to do so with varying degree of success. A basic summary for most of the studies is: process the MRI scan into a usable format, apply a feature extraction algorithm on the MRI scan to

select the appropriate features, then finally use a classification algorithm for the diagnosis of schizophrenia.

The classification of schizophrenia patients and healthy controls from sMRI scans in two large independent samples was studied in [22]. The authors used whole-brain grey matter densities from MRI scans with SVM as the classifier and concluded that SVM models trained with less than 130 samples results in an unstable model. The key difference of the study [22] from previous similar studies [23, 5], was utilizing a large dataset and using an entirely separate dataset to perform the validation. Additionally, noting that typical schizophrenia medications affect the striatum (part of the brain), they masked it out, ensuring the model doesn't relate medication effects to Schizophrenia detection. Each imaging technique provides a different view of the brain functioning. To get the benefit of different imaging techniques, a multimodal classification model [29] combined 3 different data types: resting state functional MRI (rs-fMRI), Diffusion tensor imaging (DTI) and sMRI. While the idea proposed in [29], don't have the main focus on classification but to design and evaluate a multivariate method which can find cross-information in more than two data types. In [28], multi-set canonical correlation analysis (MCCA) was used to combine the data of rs-fMRI, Electroencephalogram (EEG) and sMRI and proposed ensemble feature selection approach which resulted in very high prediction performance approaching 100% by utilizing the additional modalities. Though they also concluded that combining multiple modalities does not always result in an enhanced result. A similar study [17] which used rs-fMRI and sMRI with a similar outcome of increased accuracy when compared to single modalities. It can easily be realized what the major downfall of this multi-modal training scheme is: lack of data. Some datasets combined from various sources reach the 250 marks, like the one used in [22] or in [24], but most datasets sit at 80 samples.

Gaining insight from the previous studies about the feasibility and reliability of individual classification based on the sMRI, a novel machine-learning algorithm provided an interpretable brain signature [24]. Instead of behaving like a "Black Box" spewing out predictions, the model provided insights into the neuroanatomic markers aiding clinical interpretability. This was achieved using ElasticNet Total Variation (Enet-TV) [6] penalty, which gave Structured sparsity (which is a sparse and structured pattern of predictors). Furthermore, they pitted the Enet-TV against other SVM algorithms, showing similar predictive performance across the board. In addition to providing a clinically interpretable model, their research also suggested a shared neuroanatomical signature for early or late-stage Schizophrenia patients. The study in [17] attempted to identify the significantly contributing brain regions by averaging the weights across the five datasets used and reported the top 10 brain regions.

Until now, the papers discussed have undertaken the entire brain, sometimes obfuscating certain regions to remove the effects of medications. But studies like [18] suggest that only particular brain regions, i.e. caudate nuclei, thalami and right side amygdala, are significant in identifying a Schizophrenia patient. Since Schizophrenia patients have structural changes in hippocampus and amygdala regions, and the study [5] extracted only the hippocampus and amygdala regions of the brain for the classification of Schizophrenia subjects and concluded that hippocampal and amygdaloid structures could be utilized for classification.

3. Subjects and Methods

3.1. Subjects

The data used in this study is obtained from the Center for Biomedical Research Excellence (COBRE) data set (Available at <http://fcon.1000.projects.nitrc.org/indi/retro/cobre.html>). Data consists of 72 subjects with schizophrenia (38.1 ± 13.9 years old, range 18 – 65 years) and 74 age-matched healthy control (35.8 ± 11.5 years old, range 18 – 65 years).

3.2. 3D MRI Processing

Image processing was performed using the CAT12 package (<http://dbm.neuro.uni-jena.de>) implemented in the Statistical Parametric Mapping (SPM) toolbox version 12 (<http://www.fil.ion.ucl.ac.uk/spm/software/spm12/>). In summary, 3D T1-weighted MRI scans were parcellated into grey matter (GM), white matter (WM) and cerebrospinal fluid, skull, scalp and air cavities. In this study, the GM and WM tissues have been examined. Using a high-dimensional Diffeomorphic Anatomic Registration Through Exponentiated Lie algebra algorithm (DARTEL), the GM and WM images were normalized into Montreal Neurological Institute (MNI) space. The smoothed GM and WM images were generated through an 8-mm full-width-half-maximum Gaussian kernel. The GM and WM images were visually inspected after each step in the preprocessing phase. Besides, we used the quality check procedure implemented in the CAT12 toolbox to identify possible outliers. Finally, we resampled GM and WM images to 4-mm isotropic spatial resolution and extracted GM and WM voxel values from whole-brain data (i.e., total of 29,852 voxel values per modality) as raw features for classification tasks.

3.3. Classification Algorithms

The classification algorithms evaluated in this study for schizophrenia are explained below. Detailed information on the algorithms is available in the Supplementary file.

3.3.1. Random Forest (RaF) [1]

Random Forests proposed by Leo Breiman et al. in 2001 is a collection of tree predictors (also called tree-structured classifiers, which at their core are nested if-else statements used to vote for classes) where each tree is generated via independent and identically distributed random vectors. It has been shown that a sufficiently large forest always converge and a forest generated using random features generally produces better accuracy than a single tree classifier. RaF combined both the concepts of bagging and random subspace which improved its generalisation performance.

3.3.2. Oblique RaF (MPRaF-T, MPRaF-P and MPRaF-N) [40]

Oblique RaF was proposed to handle multiclass classification with an improved geometric property. Multisurface Proximal Support Vector Machine (MPSVM)[20] is used to generate clustering hyperplanes at the non-terminal nodes of a decision tree. Now, RaF is implemented using the MPSVM-based decision trees and then subsequently using various regularisation methods. It was shown that Oblique RaF performs better than RaF and have significantly less variance and bias. MPRaF-T, MPRaF-P, and MPRaF-N represent the MPSVM-based RaFs with Tikhonov, axis-parallel, and NULL space regularization, respectively.

3.3.3. Heterogeneous RaF [11]

As noted by [40], RaF's data splitting leads to axis-parallel decision boundaries, which can lead to poor utilization of the geometric property of the data. But [11] noted that even though Oblique RaF allows for oblique splits, it is sub-optimal. The Heterogeneous RaF uses diverse linear-classifiers at the tree's nodes and searches for the best split at every node by optimizing the impurity criteria. Heterogeneous RaF Forests are shallower and faster to train than RaF. For the decision trees, each split is rated based on impurity criterion. All the splits at each non-leaf nodes are linked with an impurity measure. The one which is having the maximum value is the selected split for that particular node. The six different classifiers which have been employed are SVM, MPSVM, Linear Discriminant Analysis (LDA), Least Squares SVM (LSSVM), Ridge Regression (RR) and Logistic Regression (LR) as they have performed well in several domains [4].

3.3.4. Kernel Ridge Regression (KRR) [10]

One of the kernel-based methods is the Kernel Ridge Regression (KRR). The KRR has a closed-form solution which lends it to faster training. Despite being relatively straightforward than other members of kernel-based methods such as SVM, it can produce comparable results. The kernel ridge regression method is based on Ridge Regression and Ordinary Least Squares.

3.3.5. K nearest neighbours (KNN) [7]

K nearest neighbours algorithm assigns the label depending upon the similarity of the point with its neighbours. A constant K is first chosen for the algorithm. The Euclidean distance of the given point is calculated and the K nearest members are selected from it. The number of data points is counted and new data points are assigned to the category for which there are maximum number of neighbours. The number of nearest neighbours, K , in our case is 5.

3.3.6. Neural Networks [13]

Neural networks are network of node layers comprising of an input layer, multiple hidden layers and an output layer. Each layer has multiple number of nodes and the nodes of each layer are interconnected with the other layers. The output of each of the layer is calculated through an activation function and the output of activation layer of the last layer is the final output. Adam optimization technique [13] has been used in order to tune the parameters.

3.3.7. Random vector functional link network (RVFL) [41]

RVFL is the randomized version of the functional link neural network. It shows that from the input layer to the hidden layer, the value of weights can be generated randomly in a suitable domain and fixed in the learning stage. The closed-form based RVFL obtains the output weights in a single-step and exhibits a higher efficiency than the iterative method.

3.3.8. Random vector functional link network with Auto Encoder (RVFLAE)[42]

Autoencoder is an unsupervised learning model for which the output and input layers share the same neurons in order to reconstruct its own inputs. In this method, we adopt a

sparse autoencoder to learn appropriate network parameters of RVFL, which are developed via l_1 norm optimization instead of the usual l_2 norm retaining more informative features.

3.3.9. Support vector machine (SVM) [2]

SVM is a binary classification algorithm which classify the labelled data in two classes. SVM generates an optimal hyperplane using data to separate the classes. Since there may be more than one hyperplane possible for that, SVM finds the optimal hyperplane to do the classification by solving a Quadratic Programming Problem (QPP). Thus, a new data point can be classified based on the optimal hyperplane formed by SVM. And if the given data is not linearly separable, SVM do the task using kernel method.

3.3.10. Twin support vector machine (TWSVM) [9]

Inspired from SVM, TWSVM is a classification algorithm which classifies the given labelled data into two classes by generating two non-parallel hyperplanes. TWSVM solves two smaller sized QPPs, unlike SVM, each for one class. The two required planes are formed by solving a problem which minimize the distance of points of corresponding class to the plane and keep it as far as possible from another class. Then, a new data point is assigned a class by calculating its distance from the planes.

3.3.11. Twin bounded support vector machine (TBSVM) [25]

TBSVM is an improved version of TWSVM which modified the optimizations problem in TWSVM to give better performance and thus making the classification more accurate. It added an extra regularization term in the formulation of TWSVM which applied the structural risk minimization principle in the model.

3.3.12. Least square twin support vector machine (LSTSVM) [15]

LSTSVM is the least squares version of TWSVM. The formulation of LSTSVM leads to fast and simple algorithm to generate the two non-parallel hyperplanes for binary classification. The two primal problems used to find the required hyperplanes are formulated in least squares sense i.e. using the equality constraints instead of the inequality constraints. The problem in LSTSVM can be solved very easily and simply by solving a system of two linear equations.

3.3.13. Robust energy based least square twin support vector machine (RELSTSVM) [31]

Robust energy based LSTSVM, proposed by Tanveer et al., is another extension of TWSVM which adds a maximum margin regularization term in primal problem and moreover uses an energy parameter in the constraints, which helped in lessening the effect of noise in the data. According to a recent study [30], RELSTSVM model leads to better classification performance among the TWSVM models.

3.3.14. Pinball general twin support vector machine (PinGTSVM) [33]

Pinball general twin support vector machine also generates non-parallel hyperplane for classification, similar to TWSVM, but uses pinball loss function in place of hinge loss without

affecting the computational complexity of the algorithm. The use of pinball loss function makes it less sensitive to noise in classification of data and make it more stable for re-sampling of data.

3.4. Feature Selection Methods

The feature selection methods evaluated in this study for schizophrenia are explained below:

3.4.1. *RankFeatures()* function and its various criterions

The *rankfeatures()* is a MATLAB[®] [21] function which ranks key features by class separability criterion. It uses various independent evaluation criteria to assess the significance of features. The criterion here refers to an objective function that minimises the overall feasible feature subset. The *rankfeatures()* uses the following feature independent criteria:

T-test [34]. The “ttest” is the *default*, criteria used by the *rankfeatures()* function. The T-test ranks the features based on an absolute value, two-sample T-test with pooled variance estimate. The “ttest” criterion assumes that the classes are normally distributed.

Entropy [14]. The “entropy” criterion uses the relative entropy, also known as Kullback-Liebler distance or divergence. The “entropy” criterion assumes that the classes are normally distributed.

Bhattacharyya [34]. The “Bhattacharyya” criterion uses the minimum attainable classification error, or Chernoff bound to rank features. The “Bhattacharyya” criterion assumes that the classes are normally distributed.

ROC [34]. The “ROC” criterion uses the area between the empirical receiver operating characteristic (ROC) curve and the random classifier slope to rank features. The “ROC” criterion is a non-parametric test.

Wilcoxon [36]. The “Wilcoxon” criterion uses the absolute value of the standardised U-statistic of a two-sample unpaired Wilcoxon test, also known as Mann-Whitney, to rank features. The “Wilcoxon” criterion is also a non-parametric test.

3.4.2. *Minimum redundancy maximum relevance (MRMR) algorithm* [3]

The MRMR algorithm is a sequential feature selection method that finds an optimal set of mutually and maximally dissimilar features. The MRMR performs this by maximising the relevance of the feature set to the response variable and minimising the redundancy of a feature set.

3.4.3. *Neighborhood Component Analysis (NCA)* [38]

NCA is a non-parametric feature selection method used explicitly for regression and classification algorithms. The feature weights (importance of a feature) are obtained using a gradient ascent technique to maximise the expected leave-one-out classification accuracy with a regularisation term.

3.5. Validation, Experimental Setup

This study used Matlab[®] R2021a[21] to implement all the required code for the different methods. The functions used were (but not limited to): *rankfeatures()*, *fscmrnr()*, *fscnca()*. In all experiments, 10-fold cross-validation was used. To study the variation of accuracy with increasing number of features and to obtain the minimum or the optimal number of features, we experimented with 100–1300 (with a step size of 100) selected features. The various hyper parameter ranges used for various methods have been tabulated in Table 1. The classification accuracies corresponding to different classification models versus feature selection approaches for the combined matter at 500 features are available in Table 2. The results of 500 features are presented as maximum accuracy for Integrated GM and WM occurs at the same. The Tables corresponding to WM and GM are available in the Supplementary file as Table S-2 and S-3.

Table 1: Parameter Ranges and values for various methods. (¹ $i = 1, 2, 3, \dots$, ²Number of Samples)

Parameter Name	Symbol	Range/Value
Penalty parameters (for TSVM-based models)	c_i^1	$\{10^i i = -5 : 5\}$
Non-linear kernel parameter	γ	$\{2^i i = -10 : 10\}$
NCA regularisation term	λ	$\{2^i / N^2 i = 1 : 20\}$
RELSTSVM parameter	E	$\{0.5, 0.6, 0.7, 0.8, 0.9, 1\}$
RVFL & RVFL-AE parameters	C N	$\{-5 : 1 : 14\}$ $\{3 : 20 : 203\}$
Ensemble size of the trees for RaF methods		100
pinGTSVM parameter	ϵ	0.05

4. Results

The performance of the classification models varies with different feature extraction methods and the number of features selected. We discuss the performance of the models with GM, WM and the integration of GM and WM data. When discussing classification models, accuracy means the average accuracy across both feature extraction techniques and the feature range, unless otherwise stated. Similarly, the accuracy assigned to a given feature extraction technique represent the average accuracy across the different classification models. Tables of other metrics i.e. AUC, Sensitivity, Specificity, Precision, F-Measure and G-Mean are available in the Supplementary file.

4.1. White Matter

Across the entire feature range (i.e. 100 - 1300), the TWSVM based classification models showed the highest average accuracy compared to the rest of the classification models, as can be seen in figure S-1a. The rest of the families, i.e. RaF, Neural Networks, KNN and KRR, stay together through the feature range. Among SVM-based classifiers, the non-linear RELSTSVM (75.16%) and non-linear TBSVM (74.30%) achieved the highest average

accuracy. Non-linear TWSVM followed them with 73.24% accuracy, closely followed by non-linear LSTSVM at 72.47%. The linear TBSVM and RELSTSVM showed $\sim 70\%$ accuracy. The lowest-performing models are pinGTSVM and RVFLAE with $\sim 60.7\%$ accuracy. Heterogeneous-RaF achieves the maximum accuracy of 84.04% for WM with 900 features selected using Wilcoxon feature selection.

Among the RaF based models, Heterogeneous-RaF shows the best performance, with an average accuracy of 67.70%. In contrast, MP RaF-T shows the lowest average accuracy with 62.93%. Among the variants of neural networks, the standard neural network and randomized based neural network show $\sim 65.5\%$ average accuracy. The autoencoder based RVFL model showed the lowest average accuracy (60.89%) among the neural network models. Also, the non-linear kernel-based KRR model (accuracy 67.15%) is better than the linear kernel-based KRR model (accuracy 63.36%) in terms of average accuracy.

Discussing the feature selection methods we can refer to S-2a, the Wilcoxon is the best choice across the entire feature range with an average accuracy of 76.16%. The NCA, in addition to being the worse performing method with an average accuracy of 61.50% it is also unstable with the number of features selected. ROC, Entropy and Bhattacharya all perform with an average accuracy of $\sim 67.1\%$. MRMR and T-test follow with 65.42% and 62.51% average accuracy, respectively.

Observing that Wilcoxon feature selection performs better than all other methods for the entire feature range, comparing classifiers based on Wilcoxon feature selection is more beneficial than average. Since we have already discussed classifiers with respect to (w.r.t.) average, we now look at classifiers which significantly different. When comparing w.r.t. Wilcoxon, the Neural Network becomes the best classifier (averaged across all feature range) at 80.91%. Heterogeneous-RaF also performs significantly better, achieving 80.42% at rank 4. Additionally, RaF-LDA also performs much better with 79.31% and rank 6.

4.2. Grey Matter

The grey matter has some exciting results both in terms of classification techniques and feature selection methods. From Figure S-1b, it is obvious to see that the SVM family of classifiers perform better than the rest of the families for the entire feature range. Another observation can be made for the neural networks that they perform worse than all other families throughout the feature range. Non-linear TBSVM and non-linear RELSTSVM achieved the highest 73.36% and 72.64% average accuracy, respectively, followed by non-linear TWSVM (72.03%) non-linear LSTSVM (70.79%). The lowest-performing models are pinGTSVM and RVFLAE with $\sim 59.5\%$ accuracy. Non-linear RELSTSVM achieves the maximum accuracy of 83.99% for grey matter with 1200 features selected using Bhattacharyya feature selection.

Among the RaF based models, the standard RaF method shows the best performance, with an average accuracy of 67.09%. In contrast, MP RaF-N shows the lowest average accuracy with 64.27%. Among the variants of neural networks, the standard neural network and randomized based neural network show $\sim 65\%$ average accuracy. The RVFLAE model showed the lowest average accuracy (59.06%) among the neural network models. Also, the

non-linear kernel-based KRR model (accuracy 69%) is better than the linear kernel-based KRR model (accuracy 63.67%) in terms of average accuracy.

Discussing the feature selection methods we can refer to S-2b, the Wilcoxon is the best choice across the entire feature range with an average accuracy of 75.21%. The NCA again is very unstable but performs very well with an average accuracy of 68.08%. Entropy and Bhattacharya both perform terribly at lower features (dipping as low as 52%) but approach the best performing Wilcoxon at higher features. MRMR performs better than both ROC and T-test at 67.5%. ROC and T-test achieve an average accuracy of 62.62% and 63.32%, respectively.

Wilcoxon feature selection performs better than all other methods for the entire feature range. Thus, we can compare classifiers based on Wilcoxon feature selection. When comparing w.r.t. Wilcoxon, the linear RELSTSVM (78.85%) perform just marginally lower than the non-linear variation (79.13%). The neural network (78.28%) performs significantly better than when using the average, but it is not at the top.

4.3. Integrated GM and WM

The combined matter (i.e. integrated GM and WM data) achieves better results than individual grey matter and white matter. From Figure 1, it can be inferred that the SVM based classifiers perform significantly better than the rest of the families for the entire feature range. The RaF based methods come next, followed by KNN and KRR. Neural networks again perform the worst. Non-linear TBSVM and non-linear RELSTSVM achieved the highest 78.47% and 77.71% average accuracy, respectively, followed by non-linear TWSVM (77.62%) non-linear LSTSVM (76.32%). The lowest-performing models are RVFLAE, pinGTSVM and KNN with ~66% accuracy. The standard neural network achieves the maximum accuracy of 86.71% for combined matter with 500 features selected using Wilcoxon feature selection.

Among the RaF based models, the Heterogeneous-RaF method shows the best performance, with an average accuracy of 73.31%. In contrast, MP RaF-T shows the lowest average accuracy with 68.92%. Among the variants of neural networks, the standard neural network and randomized based neural network show 73% and 70.77% average accuracy, respectively. The RVFLAE model showed the lowest average accuracy (65.85%) among the neural network models. Also, the non-linear kernel-based KRR model (accuracy 72.13%) is better than the linear kernel-based KRR model (accuracy 71.15%) in terms of average accuracy.

Discussing the feature selection methods we can refer to 2, the Wilcoxon is the best choice across the entire feature range with an average accuracy of 77.12%. The NCA is unstable and performs poorly with an average accuracy of 69.78%. Entropy and Bhattacharya both start out being the worst performing methods at lower features (dipping as low as 60%) but approach the best performing Wilcoxon at higher features. ROC and T-test perform better than MRMR, which is at 70.71%. ROC and T-test achieve an average accuracy of ~72%.

Wilcoxon feature selection performs better than all other methods for the entire feature range. Thus, we can compare classifiers based on Wilcoxon feature selection. When comparing w.r.t. Wilcoxon, the standard neural network comes out on top (83.98%), becoming the best classifier. The Heterogeneous-RaF becomes the second-best classifier at 82.36%.

The linear TBSVM (80.61%) and linear RELSTSVM (79.75%) performs slightly better than non-linear TBSVM (80.56%) and non-linear RELSTSVM (79.62%).

Figure 1: Average performance of classifiers families w.r.t Integrated WM and GM (Combined matter)

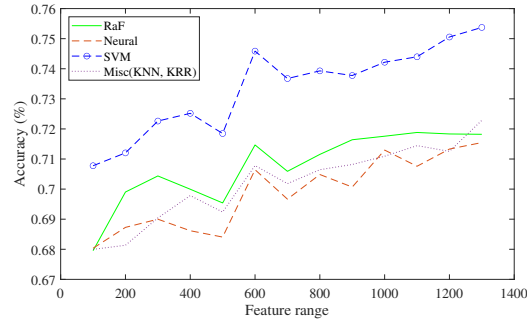


Figure 2: Average performance of feature selection methods w.r.t Integrated WM and GM (Combined matter)

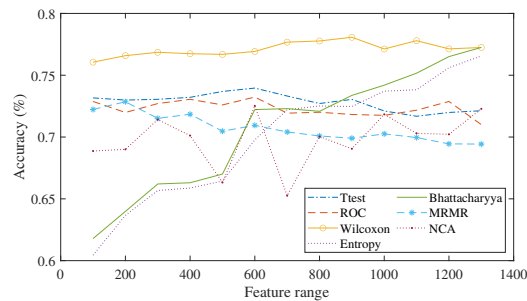
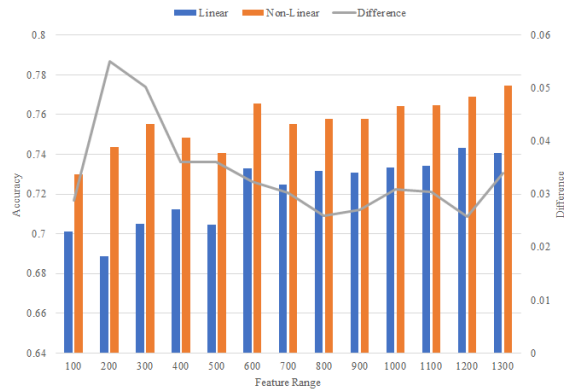


Figure 3: Performance of linear and non-linear kernels for combined matter



5. Discussion

This study aimed to perform a comprehensive evaluation of modern classification techniques and feature selection methods for schizophrenia classification. We have presented a basic overview of the different classification methods and evaluated them against different feature selection approaches on the same dataset. Selection of the optimal features and using the best available classification technique is essential for an MRI-based machine learning

Figure 4: Average performance of classification algorithms with different matter types (i.e. WM, GM and Integrated GM and WM (CM))

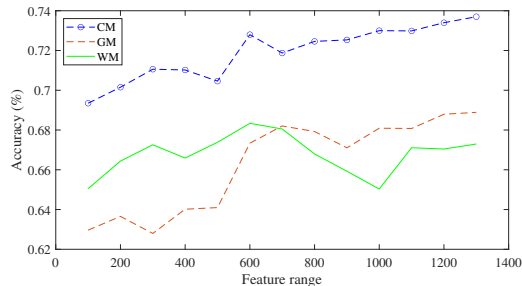


Table 2: Accuracies for Integrated GM and WM for 500 features. Feature selection methods

Classification techniques	T-Test	ROC	Wil-oxon	Entr-opy	Bhatta-charyya	MRMR	NCA
Linear variations:							
Het-RaF	72.57	73.95	82.1	67.95	62.29	74.62	63.62
KNN	69.95	71.81	73.9	60.48	60.48	65	69.67
MPRaF-N	73.1	71.67	72.48	64.33	62.95	69	61.62
MPRaF-P	71.71	75.19	75.9	61.57	67.19	70.33	69
MPRaF-T	74.52	75.86	75.1	55.29	53.33	68.29	69.1
Neural	77.38	78.76	86.71	61.57	64.24	72.29	66.52
pinGTSVM	68.57	67.76	75.1	58.95	58.86	58.1	64.14
RaF-LDA	75.33	73.29	81.38	54.62	64.33	69.71	65.81
RaF-PCA	74.57	74.48	71.62	66.29	68.33	69.81	69.9
RaF	71.67	71	75.14	64.48	65.05	75	71.05
RVFLAE	62.48	58.86	71.1	60.57	62.71	66.48	58.1
RVFL	73.33	72.76	76.48	71.19	65.57	68.48	60.95
KRR	73.14	71.86	75.24	62.76	64.1	74.43	67.1
LSTWSVM	71.85	65.62	75.4	69.28	68.25	67.33	63.71
RELSTSVM	74.28	73.18	78.84	69.12	70.7	72.14	67.8
SVM	77.43	71.95	71.71	66.24	66.9	67.1	61.71
TBSVM	75.52	73.9	80	68.29	69.05	72.95	67.86
TWSVM	73.48	67.81	78.57	68.29	68.29	69.86	66.48
Non-linear (Gaussian Kernel) variations:							
KRR	73.19	73.9	72.48	66.43	67.14	73.19	67.76
LSTWSVM	74.6	74.52	79.84	78.38	76.9	73.26	69.75
RELSTSVM	76.21	76.32	79.83	79.37	79.7	76.56	69.75
SVM	77.48	71.95	71.71	61.48	64.81	66.29	61
TBSVM	78.19	78	79.95	79.52	78.86	76.62	69.76
TWSVM	77.95	77.95	79.95	78.1	78.1	74.52	69.62

system aimed at early diagnosis. Our findings indicate that using twin SVM-based methods such as RELSTSVM or TBSVM performs best for nearly all matter types. The random forest based methods generally perform mediocly. The worst performing classification models are the pinGTSVM, RVFL-AE, KRR, KNN and MPraF-T.

Figure 3 is constructed by averaging only the methods that have both linear and non-linear variants (i.e. SVM, KRR, TWSVM, TBSVM, LSTWSVM and RELSTSVM). Inferring from figure 3, a critical observation can be made for the performance of non-linear kernel functions against linear ones. In nearly all studied methods (pinGTSVM and standard SVM being the exception), the non-linear kernel function performs better than the linear variation. This observation follows suit with the result that a linear kernel is the degenerate version of the non-linear (RBF or Gaussian) kernel [12]. Thus an adequately tuned non-linear kernel consistently out-performs the linear kernel. But, an observation can be made that with an increasing number of features, the advantage of using a non-linear kernel diminishes. This diminished performance can be attributed to the fact that at a higher number of features, one may not need to map features to a higher dimension [8]. Thus, the time penalty to tune the kernel function (in case of non-linear) becomes outweighed by the rapid computation of the linear kernel when both kernels provide relatively similar accuracy. Thus, with a large number of features, one is better off using a linear kernel.

The feature selection methods significantly impact the performance of classification models, as one might expect. The Wilcoxon was the all-around best feature selection method, performing the best for all the different matter types and across the entire feature range. This observation is supported by previous studies [35] and [19]. The exceptional performance of Wilcoxon is especially prevalent when sample sizes are small or the data doesn't resemble Normal distribution. Entropy and Bhattacharya are fascinating methods. At lower feature numbers, they perform equally terribly, but at a higher number of features, they approach the best performance. This behaviour can be seen in figure S-2, especially in integrated and grey matter. The MRMR, ROC and T-Test all perform mediocly, varying based on what matter type is used. At present, using our specific dataset, our findings indicate that we can classify schizophrenia patients with a maximum of 86.71% accuracy when using a standard neural network with 500 features from combined matter, selected using Wilcoxon. The main advantages of twin SVM based models like TWSVM, TBSVM, RELSTSVM, LSTWSVM etc. over standard SVM is that they give competitive performance in terms of accuracy and reduce the computational complexity of SVM because these models generate two non-parallel hyperplanes instead of one single hyperplane in SVM which leads to solving two smaller sized Quadratic Programming Problems (QPPs) instead of one larger QPP in SVM. The paper [30] concludes that twin SVM based models performs better than other family of classifiers.

Referring to figure 4, our results strongly suggest that using both grey matter (GM) and white matter (WM), i.e. integrated matter, leads to improved performance for the classification of schizophrenia patients. The GM performs better than WM after reaching a threshold number of features (in our case 700), but the results shoot up by a substantial margin ($\sim 4\%$) when the combined matter is used. It can be seen from figure 4 that this is the case for all the evaluated classification techniques and that there are no exceptions for this observation.

6. Conclusions

In this study, we comprehensively evaluated various classification models to identify the best available machine learning model for the classification of schizophrenia subjects. We assessed 25 classification models involving the variants of support vector machines, twin support vector machines, random forest, Kernel ridge regression and neural networks. Additionally, we evaluated 7 feature selection methods: Wilcoxon, MRMR, ROC, Entropy, T-Test, Bhattacharyya and NCA. Moreover, these evaluations were conducted on the features based on grey matter (GM), white matter (WM) and the integrated GM and WM data.

The contributions from this paper are four-fold. First, we underlined how different families of machine learning algorithms perform with the schizophrenia dataset. We found that, for the most part, the Non-linear twin SVM-based family of classifiers outperform all other classifiers. This family includes (in the order of gradually worsening performance) RELSTSVM, TBSVM, TWSVM and LSTSVM. However, the pinGTSVM is the worst-performing family member, ranking the lowest across all the classifiers. On the other hand, the KNN, KRR, MP RaF-T and non-linear SVM are the lowest performers. Most RaF-based methods occupy the middle of the spectrum. An additional observation is that, for the most part, the non-linear variant of a method outperforms the linear variation.

The second contribution is that we evaluated the performance of different feature selection methods. Our results indicate that Wilcoxon is the best performing methods with a top rank across all the matter types. Entropy and Bhattacharyya have improved performance with an increasing number of features. NCA is an unstable method, although it had a good average performance. T-Test, ROC, MRMR are also reasonable choices for feature selection, but they are not recommended.

Third, we found that utilising both grey and white matter for classification yields better results than any individual matter type.

In conclusion, the feature selection method, the number of features selected and the classification model should be appropriately chosen for better generalisation performance on the classification of schizophrenia subjects. This study recommends using standard neural network, RELSTSVM, TWSVM, or heterogeneous-RaF as the classification model with $\sim 700 - 1200$ features selected via Wilcoxon feature selection method for better generalisation performance on the classification of schizophrenia datasets. We hope that the evaluation presented in this paper encourages future research to use better classification algorithms and feature selection algorithms for clinical dataset classification.

New developments in machine learning are rapid and can improve the results of previous algorithms by a significant margin. In the future, much scope remains for the development of better specific models. Therefore, these new variants or methods need to be tested on real-life datasets such as schizophrenia to grasp their viability. In the future, one can extend this study by various margins, i.e. (1) the dataset can be enhanced to utilise data from multiple sources; thus, it should be evaluated if combining data from various sources (i.e. MRI images with varying scanning parameters such as slice thickness, field of view, bandwidth, repetition time, etc.) leads to better generalisation or if it results in a worse performance. (2) This

study utilised a single feature extraction technique (i.e. DARTEL); thus, future research for the effect of different feature extraction techniques needs to be conducted. (3) This study used a single modality (i.e. sMRI) and usage of different modalities including the Functional MRI (fMRI), Electroencephalogram (EEG) should also be evaluated using a similar setup in future studies. The source code will be available at <https://github.com/mtanveer1>

Acknowledgment

This work is supported by Science and Engineering Research Board (SERB), Government of India under Ramanujan Fellowship Scheme, Grant No. SB/S2/RJN-001/2016, and Council of Scientific & Industrial Research (CSIR), New Delhi, INDIA for funding under Extra Mural Research (EMR) Scheme grant no. 22(0751)/17/EMR-II. We gratefully acknowledge the Indian Institute of Technology Indore for providing the required facilities and support for this work.

References

- [1] L. Breiman, Random forests, *Machine Learning* 45 (1) (2001) 5–32.
- [2] C. Cortes, V. Vapnik, Support-vector networks, *Machine Learning* 20 (3) (1995) 273–297.
- [3] C. Ding, H. Peng, Minimum redundancy feature selection from microarray gene expression data, *Journal of Bioinformatics and Computational Biology* 3 (02) (2005) 185–205.
- [4] M. Fernández-Delgado, E. Cernadas, S. Barro, D. Amorim, Do we need hundreds of classifiers to solve real world classification problems?, *Journal of Machine Learning Research* 15 (1) (2014) 3133–3181.
- [5] Y. Guo, J. Qiu, W. Lu, Support vector machine-based schizophrenia classification using morphological information from amygdaloid and hippocampal subregions, *Brain Sciences* 10 (8) (2020) 562.
- [6] F. Hadj-Seleem, T. Lofstedt, E. Dohmatob, V. Frouin, M. Dubois, V. Guillemot, E. Duchesnay, Continuation of nesterov’s smoothing for regression with structured sparsity in high-dimensional neuroimaging, *IEEE Transactions on Medical Imaging* 37 (11) (2018) 2403–2413.
- [7] D.J. Hand, Principles of data mining, *Drug safety* 30 (7) (2007) 621–622.
- [8] C.W. Hsu, C.C. Chang, C.J. Lin, A practical guide to support vector classification, 2003.
- [9] Jayadeva, R. Khemchandani, S. Chandra, Twin support vector machines for pattern classification, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29 (5) (2007) 905–910.
- [10] R. Katuwal, P.N. Suganthan, An ensemble of kernel ridge regression for multi-class classification, *Procedia Computer Science* 108 (2017) 375–383.
- [11] R. Katuwal, P.N. Suganthan, L. Zhang, Heterogeneous oblique random forest, *Pattern Recognition* 99 (2020) 107078.
- [12] S.S. Keerthi, C.J. Lin, Asymptotic behaviors of support vector machines with gaussian kernel, *Neural Computation* 15 (7) (2003) 1667–1689.
- [13] D.P. Kingma, J. Ba, Adam: A method for stochastic optimization, *arXiv preprint arXiv:1412.6980* (2014).
- [14] S. Kullback, *Information theory and statistics*, Courier Corporation, 1997.
- [15] M.A. Kumar, M. Gopal, Least squares twin support vector machines for pattern classification, *Expert Systems with Applications* 36 (4) (2009) 7535–7543.
- [16] T.M. Laursen, M. Nordentoft, P.B. Mortensen, Excess early mortality in schizophrenia, *Annu Rev Clin Psychol* 10 (2014) 425–448.
- [17] D. Lei, W.H.L. Pinaya, J. Young, T. van Amelsvoort, M. Marcelis, G. Donohoe, D.O. Mothersill, A. Corvin, S. Vieira, X. Huang, S. Lui, C. Scarpazza, C. Arango, E. Bullmore, Q. Gong, P. McGuire, A. Mechelli, Integrating machine learning and multimodal neuroimaging to detect schizophrenia at the level of the individual, *Human Brain Mapping* 41 (5) (2020) 1119–1135.

- [18] X. Li, M. Black, S. Xia, C. Zhan, H.C. Bertisch, C.A. Branch, L.E. DeLisi, Subcortical structure alterations impact language processing in individuals with schizophrenia and those at high genetic risk, *Schizophrenia Research* 169 (1) (2015) 76–82.
- [19] C. Liao, S. Li, Z. Luo, Gene selection using wilcoxon rank sum test and support vector machine for cancer classification, in: *International Conference on Computational and Information Science*, volume 1, Springer, pp. 57–66.
- [20] O.L. Mangasarian, E.W. Wild, Multisurface proximal support vector machine classification via generalized eigenvalues, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 28 (1) (2005) 69–74.
- [21] MATLAB, 9.10.0.1602886 (R2021a), The MathWorks Inc., Natick, Massachusetts, 2021.
- [22] M. Nieuwenhuis, N.E. van Haren, H.E.H. Pol, W. Cahn, R.S. Kahn, H.G. Schnack, Classification of schizophrenia patients and healthy controls from structural MRI scans in two large independent samples, *NeuroImage* 61 (3) (2012) 606–612.
- [23] A. de Pierrefeu, T. Löfstedt, C. Laidi, F. Hadj-Selem, J. Bourgin, T. Hajek, F. Spaniel, M. Kolenic, P. Ciuciu, N. Hamdani, Identifying a neuroanatomical signature of schizophrenia, reproducible across sites and stages, using machine learning with structured sparsity, *Acta Psychiatrica Scandinavica* 138 (6) (2018) 571–580.
- [24] A. de Pierrefeu, T. Löfstedt, C. Laidi, F. Hadj-Selem, J. Bourgin, T. Hajek, F. Spaniel, M. Kolenic, P. Ciuciu, N. Hamdani, M. Leboyer, T. Fovet, R. Jardri, J. Houenou, E. Duchesnay, Identifying a neuroanatomical signature of schizophrenia, reproducible across sites and stages, using machine learning with structured sparsity, *Acta Psychiatrica Scandinavica* 138 (6) (2018) 571–580.
- [25] Y.H. Shao, C.H. Zhang, X.B. Wang, N.Y. Deng, Improvements on twin support vector machines, *IEEE Transactions on Neural Networks* 22 (6) (2011) 962–968.
- [26] M.E. Shenton, C.C. Dickey, M. Frumin, R.W. McCarley, A review of MRI findings in schizophrenia, *Schizophr Res* 49 (1-2) (2001) 1–52.
- [27] L. Steardo Jr, E.A. Carbone, R. de Filippis, C. Pisanu, C. Segura-Garcia, A. Squassina, P. De Fazio, L. Steardo, Application of support vector machine on fMRI data as biomarkers in schizophrenia diagnosis: A systematic review, *Frontiers in Psychiatry* 11 (2020) 588.
- [28] J. Sui, E. Castro, H. He, D. Bridwell, Y. Du, G.D. Pearlson, T. Jiang, V.D. Calhoun, Combination of FMRI-SMRI-EEG data improves discrimination of schizophrenia patients by ensemble feature selection, in: *2014 36th Annual International Conference of the IEEE Engineering in Medicine and Biology Society, IEEE, 2014*, p. 3889–3892.
- [29] J. Sui, H. He, Q. Yu, J. chen, J. Rogers, G. Pearlson, A. Mayer, J. Bustillo, J. Canive, V. Calhoun, Combination of resting state fmri, dti, and smri data to discriminate schizophrenia by N-way MCCA + jICA, *Frontiers in Human Neuroscience* 7 (2013) 235.
- [30] M. Tanveer, C. Gautam, P.N. Suganthan, Comprehensive evaluation of twin SVM based classifiers on UCI datasets, *Applied Soft Computing* 83 (2019) 105617.
- [31] M. Tanveer, M.A. Khan, S.S. Ho, Robust energy-based least squares twin support vector machines, *Applied Intelligence* 45 (1) (2016) 174–186.
- [32] M. Tanveer, T. Rajani, R. Rastogi, Y. Shao, Comprehensive review on twin support vector machines, *arXiv preprint arXiv:2105.00336* (2021).
- [33] M. Tanveer, A. Sharma, P.N. Suganthan, General twin support vector machine with pinball loss function, *Information Sciences* 494 (2019) 311–327.
- [34] S. Theodoridis, K. Koutroumbas, Chapter 5 - feature selection, in: S. Theodoridis, K. Koutroumbas (Eds.), *Pattern Recognition*, Academic Press, Boston, 2009, Fourth Edition edition, pp. 261–322.
- [35] C.j. Tian, J. Lv, X.f. Xu, Evaluation of feature selection methods for mammographic breast cancer diagnosis in a unified framework, *BioMed Research International* 2021 (2021).
- [36] F. Wilcoxon, Individual comparisons by ranking methods, in: *Breakthroughs in statistics*, Springer, 1992, pp. 196–202.
- [37] World Health Organization on Schizophrenia, World health organization on schizophrenia, <https://www.who.int/news-room/fact-sheets/detail/schizophrenia>, Last Accessed January 2022.

- [38] W. Yang, K. Wang, W. Zuo, Neighborhood component feature selection for high-dimensional data., *JCP* 7 (1) (2012) 161–168.
- [39] E. Zarogianni, T.W. Moorhead, S.M. Lawrie, Towards the identification of imaging biomarkers in schizophrenia, using multivariate pattern classification at a single-subject level, *NeuroImage: Clinical* 3 (2013) 279–289.
- [40] L. Zhang, P.N. Suganthan, Oblique decision tree ensemble via multisurface proximal support vector machine, *IEEE Transactions on Cybernetics* 45 (10) (2014) 2165–2176.
- [41] L. Zhang, P.N. Suganthan, A comprehensive evaluation of random vector functional link networks, *Information Sciences* 367 (2016) 1094–1105.
- [42] Y. Zhang, J. Wu, Z. Cai, B. Du, S.Y. Philip, An unsupervised parameter learning model for RVFL neural network, *Neural Networks* 112 (2019) 85–97.

Supplementary file of Diagnosis of Schizophrenia: A comprehensive evaluation

1 Classification Algorithms

The classification algorithms evaluated in this study for diagnosis of schizophrenia are explained below. Let the training set be $((x_1, y_1), \dots, (x_N, y_N))$ where N is the total number of training examples. X is the feature matrix $[x_1, x_2, \dots, x_N]$ of size $N \times d$ and $Y = [1, 2, \dots, m]$ is a $N \times 1$ vector of class labels.

1.1 Support Vector Machine (SVM)

Support Vector Machine [1] is supervised learning algorithm which classifies the given data into two classes by constructing a hyperplane. From the many possible hyperplanes, SVM chooses the one which maximize the margin between the two classes of data points.

SVM finds the optimal hyperplane by solving the following optimization problem,

$$\begin{aligned} \min_{w, \xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i, \\ \text{s.t.} \quad & y_i(w^T x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, m \end{aligned} \tag{1}$$

where $w, x \in \mathbb{R}^n$, ξ is the degree of misclassification and C is the penalty parameter. Using Karush-Kuhn Tucker (KKT) conditions [2], we can solve the above problem.

The optimal hyperplane is given as $w^T x + b = 0$.

SVM underperforms in places where number of features for each data point exceeds the number of training data samples.

1.2 Twin Support Vector Machine (TWSVM)

TWSVM [3] divides the given data into two classes by generating two non-parallel hyperplanes by minimizing the distance of each class points from its corresponding hyperplane. So, TWSVM solves a pair of quadratic programming problems for two classes of data points.

The pair of QPPs is given as:

$$\begin{aligned} \min_{w_1, b_1, \xi} \quad & \frac{1}{2} \|X_1 w_1 + e_1 b_1\|^2 + c_1 e_2^T \xi \\ \text{s.t.} \quad & -(X_2 w_1 + e_2 b_1) \geq e_2 - \xi, \quad \xi \geq 0 \end{aligned} \quad (2)$$

and

$$\begin{aligned} \min_{w_2, b_2, \eta} \quad & \frac{1}{2} \|X_2 w_2 + e_2 b_2\|^2 + c_2 e_1^T \eta \\ \text{s.t.} \quad & (X_1 w_2 + e_1 b_2) \geq e_1 - \eta, \quad \eta \geq 0, \end{aligned} \quad (3)$$

where c_1, c_2 are penalty parameters and ξ, η are slack variables. Solving the above primal problem by taking Lagrangian and then using K.K.T. conditions [2], we get the dual of the respective problems as follows:

$$\begin{aligned} \max_{\alpha} \quad & e_2^T \alpha - \frac{1}{2} \alpha^T G (H^T H)^{-1} G^T \alpha \\ \text{s.t.} \quad & 0 \leq \alpha \leq c_1 \end{aligned} \quad (4)$$

and

$$\begin{aligned} \max_{\gamma} \quad & e_1^T \gamma - \frac{1}{2} \gamma^T P (Q^T Q)^{-1} P^T \gamma \\ \text{s.t.} \quad & 0 \leq \gamma \leq c_2 \end{aligned} \quad (5)$$

where $H = [X_1 \ e_1], G = [X_2 \ e_2], P = [X_1 \ e_1], Q = [X_2 \ e_2]$ which can be solved to give $[w_i \ b_i], i = 1, 2$. The new data point x_j can be assigned class i by the equation $class(x_j) = \arg \min_{i=1,2} |w_i^T x_j + b_i|$

For non-linearly separable data points, appropriate kernel can be used to project data points to a higher dimensional space where they can be linearly separable.

1.3 Twin Bounded Support Vector Machine (TBSVM)

TBSVM [4] is an extension of TWSVM that introduces a maximum margin regularization term which improves the classification accuracy. The minimization of extra regularization term leads to maximizing the margin between the decision hyperplane and its parallel plane.

So, the formulation of TBSVM is given as follows:

$$\begin{aligned} \min_{w_1, b_1, \xi} \quad & \frac{1}{2} \|X_1 w_1 + e_1 b_1\|^2 + c_1 e_2^T \xi + \frac{1}{2} c_2 (\|w_1\|^2 + b_1^2) \\ \text{s.t.} \quad & -(X_2 w_1 + e_2 b_1) \geq e_2 - \xi, \quad \xi \geq 0 \end{aligned} \quad (6)$$

and

$$\begin{aligned} \min_{w_2, b_2, \eta} \quad & \frac{1}{2} \|X_2 w_2 + e_2 b_2\|^2 + c_3 e_1^T \eta + \frac{1}{2} c_4 (\|w_2\|^2 + b_2^2) \\ \text{s.t.} \quad & (X_1 w_2 + e_1 b_2) \geq e_1 - \eta, \quad \eta \geq 0, \end{aligned} \quad (7)$$

where c_1, c_2, c_3, c_4 are the penalty parameters and $\frac{1}{2}c_2(\|w_1\|^2 + b_1^2), \frac{1}{2}c_4(\|w_2\|^2 + b_2^2)$ are the extra regularization terms. Now, considering the Lagrangian and using K.K.T. conditions, one can get the dual of the above problems as:

$$\begin{aligned} \max_{\alpha} \quad & e_2^T \alpha - \frac{1}{2} \alpha^T G (H^T H + c_2 I)^{-1} G^T \alpha \\ \text{s.t.} \quad & 0 \leq \alpha \leq c_1 \end{aligned} \quad (8)$$

and for the second problem dual is

$$\begin{aligned} \max_{\gamma} \quad & e_1^T \gamma - \frac{1}{2} \gamma^T P (Q^T Q + c_4 I)^{-1} P^T \gamma \\ \text{s.t.} \quad & 0 \leq \gamma \leq c_3, \end{aligned} \quad (9)$$

where $H = [X_1 \ e_1], G = [X_2 \ e_2], P = [X_1 \ e_1], Q = [X_2 \ e_2]$ and I is the identity matrix of appropriate dimension. Solving the above dual problems one can get the required decision hyperplanes.

1.4 Least Squares Twin SVM (LSTSVM)

LSTSVM [5], proposed by Kumar et. al., is an extension of TWSVM. The formulation of LSTSVM leads to very simple and fast algorithm for generating the two non-parallel hyperplanes to classify the data points. LSTSVM

formulation solves a system of linear equations instead of two QPPs. So, the formulation of LSTSVM is given as:

$$\begin{aligned} \min_{w_1, b_1} \quad & \frac{1}{2} \|X_1 w_1 + e_1 b_1\|^2 + \frac{c_1}{2} \xi^T \xi \\ \text{s.t.} \quad & -(X_2 w_1 + e_2 b_1) = e_2 - \xi, \end{aligned} \quad (10)$$

and

$$\begin{aligned} \min_{w_2, b_2} \quad & \frac{1}{2} \|X_2 w_2 + e_2 b_2\|^2 + \frac{c_2}{2} \delta^T \delta \\ \text{s.t.} \quad & (X_1 w_2 + e_1 b_2) = e_1 - \delta, \end{aligned} \quad (11)$$

where c_1, c_2 are penalty parameters and ξ, η are slack variables. e is the vector of ones of appropriate dimension. By substituting the constraints in objective function for each problem, above equations can be solved to give a system of two linear equations as:

$$\begin{bmatrix} w_1 \\ b_1 \end{bmatrix} = -(G^T G + \frac{1}{c_1} H^T H)^{-1} G^T e_2 \quad (12)$$

$$\begin{bmatrix} w_2 \\ b_2 \end{bmatrix} = (H^T H + \frac{1}{c_2} G^T G)^{-1} H^T e_1 \quad (13)$$

where $H = [X_1 \ e_1]$ and $G = [X_2 \ e_2]$, So, the required non-parallel separating hyperplanes can be generated by solving the above system of equations.

1.5 Robust Energy-based Least Squares Twin Support Vector Machine (RELSTSVM)

RELSTSVM [6] is another extension of TWSVM which adds an extra regularization term to the objective function of each problem and incorporates the structural risk minimization principle. Moreover, RELSTSVM adds an energy to each hyperplane to lessen the impact of noise and outliers, thus making the algorithm more efficient.

Similar to LSTSVM, the constraints of RELSTSVM makes the distance of hyperplanes to the data points to be exactly 1. and the extra regularization term in objective function leads to the following formulation:

$$\begin{aligned} \min_{w_1, b_1} \quad & \frac{1}{2} \|X_1 w_1 + e_1 b_1\|^2 + \frac{c_1}{2} \xi^T \xi + \frac{c_2}{2} (\|w_1\|^2 + b_1^2) \\ \text{s.t.} \quad & -(X_2 w_1 + e_2 b_1) = E_1 - \xi \end{aligned} \quad (14)$$

and

$$\begin{aligned} \min_{w_2, b_2} \quad & \frac{1}{2} \|X_2 w_2 + e_2 b_2\|^2 + \frac{c_3}{2} \eta^T \eta + \frac{c_4}{2} (\|w_2\|^2 + b_2^2) \\ \text{s.t.} \quad & (X_1 w_2 + e_1 b_2) = E_2 - \eta \end{aligned} \quad (15)$$

where c_1, c_2, c_3, c_4 are positive parameters and E_1, E_2 are the energy parameters. Without help of any external toolbox, the above primal can be directly solved by substituting the constraint in the objective function and setting the gradient of the obtained function with respect to (w.r.t.) to w_1, b_1 equal to zero, we get the system of linear equations as

$$v_1 = -(c_1 Q^T Q + P^T P + c_2 I)^{-1} c_1 Q^T E_1 \quad (16)$$

and for the second QPP,

$$v_2 = (c_3 P^T P + Q^T Q + c_4 I)^{-1} c_3 P^T E_2 \quad (17)$$

where $v_i = \begin{bmatrix} w_i \\ b_i \end{bmatrix}$ for $i = 1, 2$, $P = [X_1 \ e]$, $Q = [X_2 \ e_2]$

1.6 Pin-GTSVM

Pin-GTSVM [7] is an extension of TWSVM that uses pinball loss instead of hinge loss which helps in making the algorithm insensitive to noise in data and more stable for re-sampling of data. Pinball loss also puts penalty on the correctly classified points. Hence, Pin-GTSVM obtains pair of planes for classification of data points by solving the following pair of QPPs:

$$\begin{aligned} \min_{w_1, b_1, \xi} \quad & \frac{1}{2} \|X_1 w_1 + e_1 b_1\|^2 + c_1 e_2^T \xi \\ \text{s.t.} \quad & -(X_2 w_1 + e_2 b_1) \geq e_2 - \xi, \\ & -(X_2 w_1 + e_2 b_1) \leq e_2 + \frac{\xi}{\tau_2} \end{aligned} \quad (18)$$

and

$$\begin{aligned} \min_{w_2, b_2, \eta} \quad & \frac{1}{2} \|X_2 w_2 + e_2 b_2\|^2 + c_2 e_1^T \eta \\ \text{s.t.} \quad & (X_1 w_2 + e_1 b_2) \geq e_1 - \eta, \\ & (X_1 w_2 + e_1 b_2) \leq e_1 + \frac{\eta}{\tau_1}, \end{aligned} \quad (19)$$

where $\tau_1, \tau_2 \in [0, 1]$ are pinball loss parameters. Forming the Lagrangian and using the K.K.T. conditions, we can obtain the dual as:

$$\begin{aligned} \max_{\alpha-\beta} \quad & e_2^T(\alpha - \beta) - \frac{1}{2}(\alpha - \beta)^T G(H^T H)^{-1} G^T(\alpha - \beta) \\ \text{s.t.} \quad & -\tau_2 c_1 e_2 \leq (\alpha - \beta) \end{aligned} \quad (20)$$

and

$$\begin{aligned} \max_{\gamma-\delta} \quad & e_1^T(\gamma - \delta) - \frac{1}{2}(\gamma - \delta)^T H(G^T G)^{-1} H^T(\gamma - \delta) \\ \text{s.t.} \quad & -\tau_1 c_2 e_1 \leq (\gamma - \delta) \end{aligned} \quad (21)$$

where $H = [X_1 \ e_1], G = [X_2 \ e_2]$ Solving the above dual one can get the two decision hyperplanes and a new data point $x \in \mathbb{R}^n$ can be assign a label for it's class using the equations of hyperplanes.

1.7 Random Forest

A random forest (RaF) is an ensemble of axis-parallel decision trees that are trained independently. Decision trees in a random forest employ recursive partitioning of the training data into smaller subsets that further aid in classification by optimizing an impurity criterion such as information gain or gini index [8]. In RaF, each non-leaf node is associated with a split function $f(x, \Theta)$ where

$$f(x, \Theta) = 1; \quad x(\Theta_1) < \Theta \quad (22)$$

$$0; \quad \textit{otherwise} \quad (23)$$

where $\Theta_1 \in \{1, 2, \dots, d\}$ is the selected feature and $\Theta_2 \in \mathbb{R}$ is a threshold. The outcome determines the child node where \mathbf{x} is routed to. The leaf nodes of the tree can either store class probability distributions or class labels based on the training samples they receive. At the time of testing, a test sample x , each tree returns probability distribution $p_t(y|x)$ and the label of the class is obtained as average or majority vote.

$$y^*(x) = \arg \max_y \frac{1}{T} \sum_{t=1}^T p_t(y|x) \quad (24)$$

Here T is the number of trees in the forest.

Random Forest requires high computational power and time as it combines numerous decision trees in order to determine the class.

1.8 Oblique Random Forest

Oblique RaF [9] was proposed to handle classification with an improved geometric property. Multisurface Proximal Support Vector Machine (MPSVM) [10] is used to generate clustering hyperplanes in decision trees. Now, RaF is implemented using the MPSVM-based decision trees and then subsequently using various regularisation methods. It was shown that Oblique RaF performs better than RaF and have significantly less variance and bias. MPRaF-T, MPRaF-P, and MPRaF-N represent the MPSVM-based RaFs with Tikhonov, axis-parallel, and NULL space regularization, respectively.

1.9 Heterogeneous Random Forest

Even though some of the oblique random forest based [9] on linear classifier perform better consistently, they aren't always the best variant of oblique random forest for every dataset. Heterogeneous RaF [11] uses several linear classifiers for generating the separating of hyperplanes. Even when some of the linear classifier based variants have lower ranks, they can still be integrated forming a heterogeneous linear classifier based oblique random forest. This would require us to evaluate n classifiers in K binary partitions, hence requiring nK number of evaluations at each node. They employed a hyper class based partitioning with one-vs-all partitioning using multiple linear classifiers at each node.

The six different classifiers which have been employed are Support Vector Machines (SVM), Multisurface Proximal SVM (MPSVM), Linear Discriminant Analysis (LDA), Least Squares SVM (LSSVM), Ridge Regression (RR) and Logistic Regression (LR) as they have performed well in several domains [12].

For the decision trees, each split is rated based on impurity criterion. All the splits at each non-leaf nodes are linked with an impurity measure Gini Index. The one which is having the maximum value is the selected split for that particular node. Instead of looking for optimal oblique split in whole search space, the recursive partitioning property exhibited by decision trees, generating few oblique splits and used their $g(i)$ for selecting the best oblique split.

The ideal gini score (g_i) is the one which is obtained when all the samples of one class are perfectly separated from other class by an oblique split. By training the linear classifiers on partitions with higher g_i and are likely to

give higher g . One can ignore the partitions with lower g_i which are likely to give lower g .

$$\min_{i=1,\dots,k} \left[\min_{i=i+1,\dots,k} \left(\frac{\text{diss}(c_i, c_j)}{\max_{m=1,\dots,k} \text{diam}(c_m)} \right) \right] \quad (25)$$

Here $\text{diss}(c_i, c_j) = \min_{x \in c_i, y \in c_j} \|x - y\|$ is the dissimilarity between cluster c_i and c_j and $\text{diam}(c) = \max_{x, y \in C} \|x - y\|$ is the intra cluster function. We use Bhattacharyya distance as the metric for distance.

1.10 Kernel Ridge Regression

Kernel Ridge Regression [13] and SVM [1] are the best known members using kernel method. Kernel based methods are very useful when there is non-linear structure in data. KRR is faster to train and simpler with its closed form solution and can achieve performance which is comparable to complex methods such as SVM.

The kernel ridge regression method is based on Ridge Regression [14] and Ordinary Least Squares. The OLS minimizes the loss $\min_{\beta} \|Y - X\beta\|^2$ which is the L_2 norm. A shrinkage parameter λ is added to control the trade-off between variance and bias in the above expression giving us the following problem.

$$\min_{\beta} \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \quad (26)$$

The closed form solution for above can be problem given as $\beta = (X^T X + \lambda I)^{-1} X^T Y$. The label predicted for the new unlabeled example x is given as $\beta^T x$. The Kernel ridge regression method extends linear regression into non-linear and high-dimensional space. The data which is present in X is replaced with the feature vectors $:x_i \rightarrow \phi = \phi(x_i)$ induced by the kernel where $K_{ij} = k(x_i, x_j) = \phi(x_i)\phi(x_j)$. Hence the new predicted class label for the new example x is given as :

$$Y^T (K + \lambda I)^{-1} k \quad (27)$$

Here $k = (k_1, k_2, \dots, k_N)^T$, $k_n = x_n \cdot x$ and $n = 1, 2, \dots, N$.

1.11 Random vector functional link network (RVFL) [15]

RVFL [15] is the randomized version of the functional link neural network. Here, the input layer to the hidden layer, the weights are generated randomly in a suitable domain and fixed in the learning stage. Weights are generated in this manner ensuring that the activation functions $g(a_j^T x + b_j)$ are not all saturated. All weights are generated with uniform distribution within $[-S, +S]$. Here S is a scale factor which is determined at the stage of parameter tuning. Only the output weights need to be determined by solving the problem :

$$y_i = d_i^T \beta, \quad i = 1, 2, \dots, N \quad (28)$$

Here P is the number of data samples, t is the target and d is the vectorised concatenation of random and the original features. Least squares can be used as a regularization technique in order to avoid over-fitting and obtain the solution. The two classes of RVFL algorithm are iterative RVFL, which obtains the output weights in an iterative manner based on the gradient of the error function and the closed-form based RVFL, which obtains the output weights in a single-step. The closed-form based RVFL exhibits a higher efficiency. L2 norm regularized least square is used to solve the following problem :

$$\sum_i (y_i - d_i^T \beta)^2 + \lambda \|\beta\|^2 ; \quad i = 1, 2, \dots, N \quad (29)$$

The solution for the same is given as $\beta = D(D^T D + \lambda I)^{-1} Y$, where λ is the regularization parameter to be tuned. D and Y are the stacked features and targets of all the data samples in matrix form.

1.12 Random vector functional link network with AutoEncoder (RVFLAE) [16]

Autoencoder is an unsupervised learning model for which the output and input layers share the same neurons in order to reconstruct its own inputs instead of predicting target values for given input data. Sparse pre-trained RVFL is the unsupervised parameter learning method for RVFL. A sparse

autoencoder is used to learn appropriate network parameters, which are developed via l_1 norm optimization instead of the usual l_2 norm. This means that more informative features will be retained to participate in the subsequent learning processes. During the learning process, the sparse autoencoder captures the excellent features in the encoding stage and learns the output weights in the decoding stage. Let the input data be X , then the sparse autoencoder has optimization problem given as :

$$O_w = \arg \min \{ \|\tilde{H}\bar{\omega} - X\|^2 + \|\bar{\omega}\|_{l_1} \}$$

Here $\tilde{H} \in \mathbb{R}^{N \times L}$ is the output matrix of hidden layer obtained via random feature mapping. $\bar{\omega} \in \mathbb{R}^{L \times d}$ is the output weight matrix of the sparse encoder. $\|\tilde{H}\bar{\omega} - X\|^2$ measures the loss to model the reconstruction process of input data and $\|\bar{\omega}\|_{l_1}$ is the l_1 norm regularization. The solution of this optimization problem is given by Fast iterative shrinkage threshold algorithm (FISTA) [17].

1.13 K nearest neighbours

K nearest neighbours algorithm assigns the label depending upon the similarity of the point with its neighbours. A constant K is first chosen for the algorithm. The Euclidean distance of the given point is calculated and the K nearest members are selected from it. The number of data points is counted and assign data points is assigned to the category for which there are maximum number of neighbours.

1.14 Neural Network

Neural networks [18] are network of node layers comprising of an input layer, multiple hidden layers and an output layer. Each layer has multiple number of nodes and the nodes of each layer are interconnected with the other layers. Each of the nodes consists of weights which are tuned by training on the examples. The output of each of the layer is calculated through an activation function, which is then passed to the next layer. We have divided the dataset into 85% training and 15% test set for our network. The layers of our network depending upon the activation functions are Feature Input Layer \rightarrow Fully Connected Layer \rightarrow Batch Normalized Layer \rightarrow Relu Layer \rightarrow Fully Connected Layer \rightarrow SoftMaxLayer \rightarrow Classification Layer.

For the tuning of our parameters, we have used Adam Optimizer. It is the combination of the Stochastic Gradient Descent with momentum and the Root Mean Square Propagation and hence is quite efficient.

The network once trained, is used to make predictions on the test set.

2 Statistical Analysis

We follow Friedman Test [19, 20] to test the significant difference among the linear and non-linear classification models. Consider k algorithms are evaluated N datasets/feature selection techniques. Let r_i^j be the rank of j^{th} algorithm on i^{th} feature selection technique. The ranks of the algorithms are compared by taking their average performance $R_j = \frac{1}{N} \sum_i r_i^j$. As per the null hypotheses, the ranks R_j of the algorithms should be equal considering that all algorithms are equivalent.

$$\chi_F^2 = \frac{12N}{(k)(k+1)} \left[\sum_j R_j^2 - \frac{k(k+1)^2}{4} \right]$$

Here the Friedman statistic is distributed as per χ_F^2 with degrees of freedom as $k-1$ considering k and N are large ($N > 10$ and $k > 5$).

Friedman's χ_F^2 is considered to be conservative and a better statistic has been derived [21]. The better statistic F_F is given as follows

$$F_F = \frac{(N-1)\chi_F^2}{N(k-1) - \chi_F^2}$$

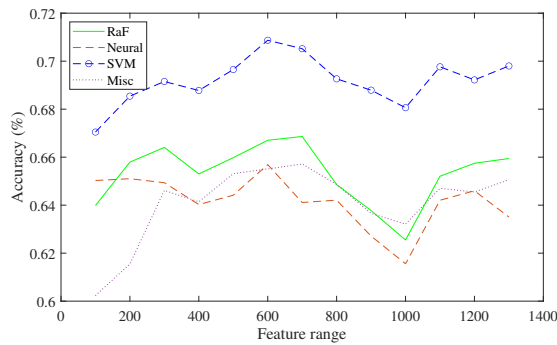
This statistic having $k-1$ and $(k-1)(N-1)$ degrees of freedom is distributed as per the F distribution.

We can proceed with a post-hoc test in case the null-hypothesis is rejected. The Nemenyi test [22] is used while comparing the classifiers with each other. If the average ranks of two classifiers differ by atleast the critical difference,

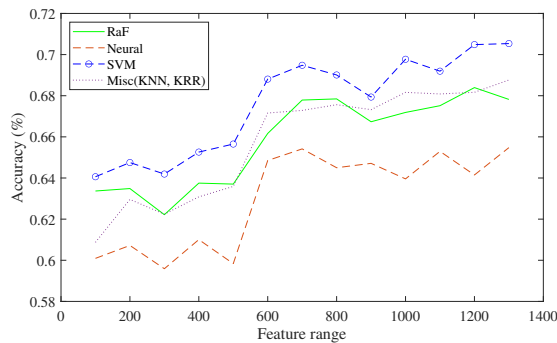
$$\text{Critical Difference} = q_\alpha \sqrt{\frac{k(k+1)}{6N}} \quad (30)$$

then the performance of classifiers is considered to be different significantly. Here, the q_α critical values are dependent upon Studentized range statistic divided by $\sqrt{2}$. Here, $N = 7, K = 12$. With simple calculation, we have

$\chi_F^2 = 54.48, F_F = 14.52$. From Statistical tables, $F_F(11, 66) = 1.937$. Since $14.52 > 1.937$, hence we reject the null hypothesis. Thus, significant difference exists among the models. With $q_{\alpha=0.05} = 3.268$, we get $CD = 6.3$. Hence, if the rank difference of two classifiers is atleast CD , then the two models are significantly different. Based on the Nemenyi test, Table S-1 gives the significant difference among the models. No significant difference exists among the other methods not given in Table S-1.



(a) White Matter

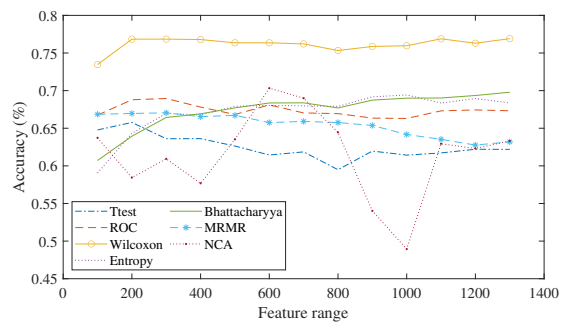


(b) Grey Matter

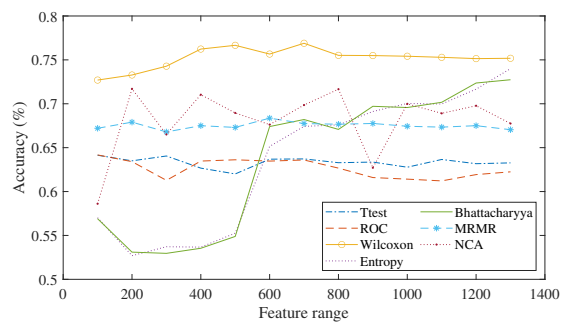
Figure S-1: Average performance of classifiers families

Table S-1: Statistical analysis on linear and non-linear variants of classification models.

Models	KRR (Linear)	KRR (Non-Linear)	LSTWSVM (Linear)	TWSVM (Linear)	SVM (Linear)	SVM (Non-Linear)
RELSTWSVM (Non-Linear)	Yes	No	Yes	No	Yes	Yes
TBSVM (Non-Linear)	Yes	Yes	Yes	Yes	Yes	Yes
TWSVM (Non-Linear)	Yes	No	Yes	No	Yes	Yes



(a) White Matter



(b) Grey Matter

Figure S-2: Average performance of feature selection methods

3 Supplementary Tables

Tables S-2 and S-3 are accuracies for White Matter (900 features) and Grey Matter (1200 features). Tables S-4, S-5, S-6, S-7, S-8, S-9 are other performance metrics (AUCs, Sensitivity, Specificity, Precision, F-Measure and G-Mean) for White Matter (900 features). Tables S-10, S-11, S-12, S-13, S-14, S-15 are other performance metrics for integrated GM and WM (500 features). Tables S-16, S-17, S-18, S-19, S-20, S-21 are other performance metrics for Grey Matter (1200 features).

References

- [1] C. Cortes, V. Vapnik, Support-vector networks, *Machine Learning* 20 (3) (1995) 273–297.
- [2] R. Fletcher, *Practical methods of optimization*, John Wiley & Sons, 2013.
- [3] Jayadeva, R. Khemchandani, S. Chandra, Twin support vector machines for pattern classification, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29 (5) (2007) 905–910.
- [4] Y.-H. Shao, C.-H. Zhang, X.-B. Wang, N.-Y. Deng, Improvements on twin support vector machines, *IEEE Transactions on Neural Networks* 22 (6) (2011) 962–968.
- [5] M. A. Kumar, M. Gopal, Least squares twin support vector machines for pattern classification, *Expert Systems with Applications* 36 (4) (2009) 7535–7543.
- [6] M. Tanveer, M. A. Khan, S.-S. Ho, Robust energy-based least squares twin support vector machines, *Applied Intelligence* 45 (1) (2016) 174–186.
- [7] M. Tanveer, A. Sharma, P. N. Suganthan, General twin support vector machine with pinball loss function, *Information Sciences* 494 (2019) 311–327.
- [8] A. Criminisi, J. Shotton, E. Konukoglu, Decision forests: A unified framework for classification, regression, density estimation, manifold

Table S-2: Accuracies for White Matter for 900 features.

Classification techniques	Feature selection methods						
	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Linear variations:							
Het-RaF	67.95	67.24	84.05	66.57	67.76	66.33	51.9
KNN	54.86	65.14	74.52	57	55.62	60.95	44.48
MPRaF-N	59.62	66.52	70.24	68.48	61.67	65	47.24
MPRaF-P	65.19	65.62	69.67	69.24	65	66.62	50.76
MPRaF-T	63.05	51.43	73.71	65.14	66.33	55.43	49.33
Neural	55.57	66.62	81.9	63.57	60.86	62.38	47
pinGTSVM	54.1	58.19	73.19	71.24	68.67	50.05	49.81
RaF-LDA	56.86	63.67	76.57	62.05	67.81	67.86	48.57
RaF-PCA	62.24	68.57	72.38	67.71	64.48	63.14	47.24
RaF	61.76	65	70.52	69.14	63	68.57	50.76
RVFLAE	50.1	58.71	72.48	60.81	58.71	60.19	56.71
RVFL	61.67	61	77.95	64.14	72.52	68.33	55.71
KRR	66.52	69.86	73.76	71	71	52	54.19
LSTWSVM	62.92	63.79	75.49	67.5	66.17	65.92	56.99
RELSTWSVM	65.88	72.63	77.79	74.06	74.06	70.01	59.54
SVM	58.86	64.38	74.43	67.57	68.24	64.29	56.67
TBSVM	66.33	72	77.29	72.43	72.38	69.71	60.71
TWSVM	59.62	65.1	77.24	68.24	68.29	66.95	59.57
Non-linear (Gaussian Kernel) variations:							
KRR	64.48	66.43	73.81	69.9	69.9	68.52	52.81
LSTWSVM	66.91	71.08	79.29	77.63	79.4	71.56	58.4
RELSTWSVM	69.25	76.15	81.09	81.11	82.13	72.68	59.89
SVM	57.62	67.14	75.1	67.48	68.24	66.48	56.62
TBSVM	68.62	73.33	79.9	79.43	78.81	73.29	61.76
TWSVM	67.24	72.71	78.52	78.71	78.71	72.67	59.71

Table S-3: Accuracies for Grey Matter for 1200 features.

Classification techniques	Feature selection methods						
	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Linear variations:							
Het-RaF	62.86	61.52	76.71	71.29	73.95	70.43	67.71
KNN	62.24	60.86	73.29	63	64.33	67.76	65.14
MPRaF-N	61.52	63.57	72.43	69.62	74.43	67.62	71.24
MPRaF-P	65	63.57	69.71	71.86	71	68.95	69.14
MPRaF-T	67.76	64.9	75.9	69.52	70.24	64.29	71.81
Neural	62.24	60.81	74.62	61.57	65	61.67	66.38
pinGTSVM	56.76	53.29	71.76	67.05	69.86	62.14	66.38
RaF-LDA	65.81	56.71	77.9	69.9	65	67.71	66.19
RaF-PCA	67	58.62	77.24	68.95	73.76	67.81	61.71
RaF	66.29	63.62	71	70.29	71.14	68.29	67.81
RVFLAE	54.86	53.43	69.71	58.38	61.76	60.38	64.29
RVFL	61.67	65.48	71.81	72.52	67	68.52	65
KRR	60.86	58.76	76.67	70.52	70.48	67.67	70.38
LSTWSVM	59.5	58.25	75.26	70.3	69.89	67.94	69.56
RELSTWSVM	64.9	64.26	78.26	73.73	73.33	69.75	73.91
SVM	57.62	52.76	74.57	71.86	73.24	67.1	71.14
TBSVM	64.24	64.29	79.48	72.48	73.9	69.29	71.81
TWSVM	60.24	58.95	74.57	72.48	67.86	67.05	66.48
Non-linear (Gaussian Kernel) variations:							
KRR	68.52	69.9	76.05	73.95	73.19	70.29	67.62
LSTWSVM	64.52	64.91	78.16	80.9	82.51	71.12	76.29
RELSTWSVM	67.95	68.78	77.88	82.82	84	69.91	77.28
SVM	55.52	60.86	74.57	72.52	75.14	60.38	73.86
TBSVM	71.19	69.86	77.95	82.24	82.86	72.43	76.71
TWSVM	67.1	68.38	77.9	82.14	82.86	71.62	76.71

Table S-4: Evaluation of classification models based on AUCs with White Matter (900 features).

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	67.98	67.29	83.68	68.51	69.14	67.58	51.23
KNN	54.06	65.40	74.88	56.24	54.97	61.32	45.83
KRR (Linear)	65.89	70.42	73.66	72.47	72.47	55.03	54.55
KRR (Non-Linear)	65.33	67.42	73.99	71.09	71.16	68.72	53.91
LSTWSVM (Linear)	57.12	59.69	68.22	64.65	63.65	57.70	55.08
LSTWSVM (Non-Linear)	60.46	68.59	88.54	81.65	78.00	65.92	51.61
MPRaF-N	63.58	67.71	73.03	70.68	64.32	66.26	51.37
MPRaF-P	66.22	65.37	69.60	70.41	66.44	67.62	51.73
MPRaF-T	63.24	56.94	74.54	66.13	67.04	57.01	49.59
Neural	56.04	66.26	81.91	64.33	62.41	61.68	46.55
pinGTSVM	55.79	59.00	73.80	69.54	67.77	50.10	48.29
RaF-LDA	57.84	64.84	76.26	62.86	67.88	67.21	49.74
RaF-PCA	62.54	68.72	72.11	67.53	64.95	62.84	47.04
RaF	62.23	64.18	71.27	69.59	63.84	69.42	51.38
RELSTSVM (Linear)	59.79	73.62	73.81	75.69	75.69	68.77	61.61
RELSTSVM (Non-Linear)	70.73	84.62	74.77	83.15	86.76	66.92	57.83
RVFLAE	49.95	57.97	72.53	61.12	58.90	60.37	56.45
RVFL	61.73	60.43	77.49	64.89	72.27	68.20	55.80
SVM	56.65	67.37	75.33	68.42	69.20	65.88	55.81
TBSVM (Linear)	65.08	71.91	76.81	73.97	73.70	72.53	60.11
TBSVM (Non-Linear)	68.81	73.44	79.71	79.27	79.65	73.38	59.92
TWSVM (Linear)	58.92	63.72	76.81	68.21	68.12	66.41	58.29
TWSVM (Non-Linear)	67.79	72.27	78.10	78.16	77.83	72.73	57.94

Table S-5: Sensitivities of the classification models for White Matter (900 features).

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	63.27	63.56	80.56	59.99	59.88	60.84	38.35
KNN	25.35	52.64	62.98	17.20	16.09	41.32	18.92
KRR (Linear)	65.13	69.07	72.19	70.23	70.23	67.75	52.38
KRR (Non-Linear)	53.96	70.05	61.21	61.12	60.00	62.29	30.03
LSTWSVM (Linear)	68.73	67.88	82.75	70.35	68.69	74.15	58.91
LSTWSVM (Non-Linear)	73.35	73.56	70.04	73.62	80.81	77.19	65.20
MPRaF-N	76.25	68.80	82.48	63.42	57.25	65.19	67.48
MPRaF-P	60.92	58.28	65.00	65.49	59.00	61.60	32.04
MPRaF-T	54.13	56.42	67.75	47.65	49.66	30.17	39.51
Neural	48.56	60.27	78.50	61.18	58.92	57.98	46.44
pinGTSVM	61.40	63.77	66.03	74.71	71.04	62.53	56.38
RaF-LDA	51.10	55.60	63.11	52.36	53.24	56.57	25.79
RaF-PCA	57.64	60.01	67.21	58.27	52.15	48.35	33.14
RaF	55.28	55.32	66.72	66.83	58.56	62.42	35.99
RELSTSV (Linear)	71.96	71.63	81.77	72.44	72.44	71.25	57.46
RELSTSV (Non-Linear)	67.77	67.67	87.40	79.06	77.51	78.44	61.95
RVFLAE	42.71	52.24	66.21	58.33	58.33	55.16	56.88
RVFL	49.22	56.24	76.21	64.42	68.61	63.32	52.27
SVM	51.34	68.41	74.85	66.92	67.94	59.18	51.46
TBSVM (Linear)	72.09	72.19	70.72	74.26	72.60	72.75	57.50
TBSVM (Non-Linear)	71.34	70.08	77.46	74.00	80.23	64.91	46.54
TWSVM (Linear)	54.09	57.61	71.72	69.06	68.88	61.73	59.11
TWSVM (Non-Linear)	73.65	82.16	73.12	65.04	64.37	64.50	44.11

Table S-6: Specificity of classification models for White Matter (900 features).

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	72.70	71.02	86.80	77.04	78.40	74.33	64.11
KNN	82.77	78.16	86.77	95.28	93.85	81.33	72.73
KRR (Linear)	66.65	71.77	75.14	74.70	74.70	42.32	56.73
KRR (Non-Linear)	76.69	64.79	86.78	81.06	82.31	75.14	77.80
LSTWSVM (Linear)	67.19	65.39	77.44	66.78	65.42	71.14	57.96
LSTWSVM (Non-Linear)	72.38	72.41	73.68	74.09	79.13	71.69	57.39
MPRaF-N	50.90	66.62	63.59	77.95	71.38	67.33	35.25
MPRaF-P	71.52	72.47	74.20	75.33	73.88	73.63	71.42
MPRaF-T	72.35	57.46	81.33	84.62	84.42	83.85	59.66
Neural	63.53	72.25	85.31	67.48	65.89	65.39	46.66
pinGTSVM	50.17	54.22	81.57	64.38	64.49	37.67	40.19
RaF-LDA	64.58	74.08	89.40	73.36	82.52	77.85	73.69
RaF-PCA	67.44	77.42	77.01	76.79	77.74	77.33	60.94
RaF	69.17	73.04	75.82	72.35	69.12	76.42	66.77
RELSTSVM (Linear)	69.89	71.27	78.65	71.04	71.04	68.72	57.58
RELSTSVM (Non-Linear)	68.33	70.59	81.81	77.86	77.53	72.94	58.82
RVFLAE	57.19	63.70	78.85	63.91	59.46	65.58	56.02
RVFL	74.24	64.63	78.77	65.37	75.94	73.09	59.33
SVM	61.96	66.33	75.81	69.91	70.46	72.59	60.15
TBSVM (Linear)	58.07	71.63	82.91	73.69	74.80	72.31	62.71
TBSVM (Non-Linear)	66.29	76.80	81.96	84.54	79.06	81.84	73.31
TWSVM (Linear)	63.75	69.83	81.91	67.35	67.35	71.08	57.46
TWSVM (Non-Linear)	61.93	62.37	83.07	91.28	91.28	80.96	71.78

Table S-7: Precision of classification models for White Matter (900 features).

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	69.45	69.61	84.71	74.13	73.58	69.92	50.08
KNN	52.83	70.18	81.75	NaN	NaN	67.98	32.95
KRR (Linear)	69.73	70.97	72.06	71.31	71.31	55.01	53.35
KRR (Non-Linear)	71.87	67.75	79.48	77.53	78.89	69.09	64.19
LSTWSVM (Linear)	59.52	60.49	71.45	64.78	63.67	61.82	54.45
LSTWSVM (Non-Linear)	63.24	68.73	79.42	76.75	75.78	67.13	52.89
MPRaF-N	63.51	68.59	67.60	75.48	69.35	66.65	50.71
MPRaF-P	68.79	68.29	69.28	74.44	70.58	72.45	52.88
MPRaF-T	68.22	NaN	75.31	74.08	79.95	72.33	50.02
Neural	56.87	68.40	82.23	63.73	60.33	63.04	44.84
pinGTSVM	57.18	56.86	75.65	68.18	66.18	49.46	47.90
RaF-LDA	59.49	70.65	86.13	68.46	76.67	71.81	NaN
RaF-PCA	64.83	72.07	74.84	71.59	71.76	66.57	44.83
RaF	64.19	70.36	72.54	70.60	67.15	74.11	50.35
RELSTSVM (Linear)	62.51	71.18	75.46	71.08	71.08	67.03	58.19
RELSTSVM (Non-Linear)	68.12	75.99	77.62	79.56	81.14	68.74	56.86
RVFLAE	43.47	58.43	74.63	60.99	58.27	61.37	57.48
RVFL	63.42	61.00	75.61	62.65	72.89	69.18	50.61
SVM	55.09	66.95	72.99	66.70	66.97	70.02	57.20
TBSVM (Linear)	65.67	70.82	79.75	71.38	71.20	71.35	59.65
TBSVM (Non-Linear)	67.36	75.31	78.52	81.40	77.62	75.26	56.96
TWSVM (Linear)	59.28	65.49	77.73	67.52	67.83	66.84	56.58
TWSVM (Non-Linear)	67.63	69.73	78.39	84.05	82.38	75.22	55.89

Table S-8: F-Measures for White Matter for 900 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	64.33	64.11	81.89	61.47	62.45	62.39	NaN
KNN	NaN	58.40	69.46	NaN	NaN	47.93	NaN
KRR (Linear)	65.47	68.22	71.02	68.54	68.54	57.01	51.34
KRR (Non-Linear)	58.46	66.42	68.33	65.49	65.29	64.73	37.97
LSTWSVM (Linear)	60.79	61.48	72.13	65.25	64.10	63.07	55.46
LSTWSVM (Non-Linear)	64.74	69.60	80.26	77.30	77.14	67.93	53.68
MPRaF-N	63.96	66.09	71.52	65.35	56.24	63.49	54.44
MPRaF-P	61.84	61.11	66.01	67.39	60.88	62.29	NaN
MPRaF-T	58.15	NaN	69.73	54.76	58.31	38.30	41.98
Neural	49.86	62.74	79.43	61.22	57.63	59.67	44.49
pinGTSVM	55.21	58.32	68.97	71.02	68.13	53.62	50.50
RaF-LDA	53.25	58.66	71.15	57.34	61.52	62.24	NaN
RaF-PCA	58.97	63.26	69.25	62.86	57.12	53.53	36.71
RaF	57.31	59.80	67.83	67.79	58.61	64.53	39.84
RELSTSVM (Linear)	63.63	71.80	75.85	72.19	72.19	67.87	58.88
RELSTSVM (Non-Linear)	68.82	76.79	77.95	80.03	81.63	69.32	57.58
RVFLAE	NaN	53.97	69.21	58.51	57.44	56.56	55.49
RVFL	53.80	58.03	75.00	61.45	69.95	65.02	NaN
SVM	51.72	66.28	72.76	65.28	65.84	60.93	52.46
TBSVM (Linear)	67.18	70.37	74.12	70.64	70.01	68.48	56.52
TBSVM (Non-Linear)	67.03	71.14	77.61	76.66	77.99	68.22	NaN
TWSVM (Linear)	55.35	59.41	74.02	66.80	67.04	63.04	56.64
TWSVM (Non-Linear)	68.53	72.79	75.35	72.09	70.98	67.96	NaN

Table S-9: G-Means for White Matter for 900 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	65.31	65.26	82.26	63.72	64.35	63.79	43.29
KNN	35.86	59.82	70.88	NaN	NaN	50.71	22.64
KRR (Linear)	66.43	69.10	71.57	69.62	69.62	59.09	52.08
KRR (Non-Linear)	60.47	67.62	69.32	67.29	67.24	65.21	41.85
LSTWSVM (Linear)	5.61	5.57	5.56	5.52	5.75	5.72	5.51
LSTWSVM (Non-Linear)	0.31	0.29	0.31	0.31	0.31	0.31	0.31
MPRaF-N	66.75	67.35	73.24	67.27	59.14	64.66	56.64
MPRaF-P	63.26	62.12	66.56	68.63	62.63	64.15	39.65
MPRaF-T	59.57	NaN	70.60	57.27	61.22	43.47	43.32
Neural	51.20	63.51	79.90	61.83	58.61	60.08	45.05
pinGTSVM	57.05	59.29	69.89	71.23	68.37	54.78	51.31
RaF-LDA	54.24	60.69	72.82	58.83	63.20	63.20	NaN
RaF-PCA	60.08	64.61	70.11	63.86	59.42	55.29	37.80
RaF	58.44	61.23	68.69	68.25	60.24	66.20	41.38
RELSTSVM (Linear)	5.50	5.38	5.40	5.36	5.44	5.42	5.54
RELSTSVM (Non-Linear)	0.35	0.33	0.33	0.35	0.34	0.34	0.34
RVFLAE	41.79	54.64	69.81	59.08	57.87	57.40	56.32
RVFL	55.01	58.32	75.45	62.47	70.34	65.63	50.69
SVM	52.45	66.97	73.33	66.03	66.63	62.71	53.38
TBSVM (Linear)	68.02	70.93	74.67	71.70	70.92	70.19	57.52
TBSVM (Non-Linear)	68.15	71.90	77.80	77.17	78.46	69.12	51.11
TWSVM (Linear)	56.00	60.44	74.37	67.54	67.69	63.65	57.23
TWSVM (Non-Linear)	69.56	74.33	75.55	73.29	72.13	68.87	48.90

Table S-10: AUCs for Combined Matter for 500 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	72.16	73.21	81.50	69.35	63.93	73.62	64.08
KNN	68.77	71.92	72.53	59.92	59.92	65.27	69.26
KRR (Linear)	72.22	71.00	74.07	63.64	64.98	73.74	66.33
KRR (Non-Linear)	71.89	73.68	72.18	66.68	67.73	73.22	67.17
LSTWSVM (Linear)	68.44	64.33	68.77	65.21	63.21	66.38	66.04
LSTWSVM (Non-Linear)	67.99	64.88	90.13	72.05	70.21	69.51	64.70
MPRaF-N	73.17	72.11	73.35	65.75	63.99	67.87	63.09
MPRaF-P	70.30	75.02	76.53	61.49	67.19	70.72	68.64
MPRaF-T	73.02	74.14	74.70	54.97	55.62	68.27	68.94
Neural	76.10	77.60	85.88	61.74	64.36	72.27	66.15
pinGTSVM	69.25	67.13	75.92	57.73	56.75	59.33	64.31
RaF-LDA	74.30	72.01	81.20	56.27	65.80	69.06	66.08
RaF-PCA	73.45	73.91	71.66	66.45	68.29	70.77	69.32
RaF	71.46	70.54	75.91	64.72	65.26	75.23	70.51
RELSTSVM (Linear)	76.63	75.20	73.35	64.32	64.01	74.20	63.14
RELSTSVM (Non-Linear)	74.96	72.52	76.87	77.98	78.65	79.55	64.70
RVFLAE	61.30	57.28	70.83	61.48	63.33	66.33	58.55
RVFL	72.02	71.68	75.85	71.72	66.38	68.34	60.06
SVM	75.92	70.34	71.13	61.65	65.11	66.41	60.80
TBSVM (Linear)	74.49	72.89	78.84	69.80	70.68	72.06	66.97
TBSVM (Non-Linear)	76.44	76.79	78.36	79.25	78.56	75.86	69.46
TWSVM (Linear)	71.51	67.15	77.92	69.80	69.58	70.23	66.36
TWSVM (Non-Linear)	76.21	76.32	77.90	77.37	77.30	73.90	68.82

Table S-11: Sensitivities for Combined Matter for 500 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	61.14	64.89	80.55	58.04	53.74	70.66	52.66
KNN	58.22	67.94	67.56	26.00	26.00	53.73	56.13
KRR (Linear)	71.35	65.93	70.08	59.34	60.77	71.70	62.32
KRR (Non-Linear)	66.47	65.38	55.99	50.60	54.95	65.70	56.52
LSTWSVM (Linear)	75.26	66.90	82.04	73.35	73.29	68.29	61.38
LSTWSVM (Non-Linear)	81.21	84.16	69.55	84.71	83.60	77.01	74.79
MPRaF-N	77.04	75.87	74.59	47.13	48.89	63.74	60.28
MPRaF-P	67.18	70.84	73.27	52.85	61.58	67.48	58.18
MPRaF-T	65.50	71.29	70.59	20.69	23.71	61.21	59.03
Neural	70.47	72.88	84.56	55.93	60.22	71.01	64.06
pinGTSVM	73.75	70.35	80.21	74.84	72.16	64.29	65.77
RaF-LDA	65.42	65.93	75.80	44.52	56.79	67.44	55.02
RaF-PCA	66.07	72.91	68.73	57.36	58.57	67.88	65.32
RaF	70.51	72.19	74.70	55.57	55.63	71.66	64.59
RELSTSVM (Linear)	71.94	71.16	84.34	73.92	77.40	70.08	72.46
RELSTSVM (Non-Linear)	77.45	80.12	82.78	80.76	80.76	73.56	74.79
RVFLAE	53.40	51.14	66.52	59.33	60.40	62.22	58.57
RVFL	66.80	62.32	75.80	64.62	60.65	67.04	57.45
SVM	70.68	63.85	69.18	52.37	65.02	65.64	59.19
TBSVM (Linear)	68.01	64.26	76.01	65.89	61.71	71.45	61.71
TBSVM (Non-Linear)	74.77	65.93	77.87	73.96	70.04	69.27	51.47
TWSVM (Linear)	66.23	61.15	75.62	65.89	63.65	69.71	65.53
TWSVM (Non-Linear)	74.96	70.86	73.35	70.21	70.21	65.58	63.67

Table S-12: Specificity's for Combined Matter for 500 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	83.19	81.53	82.45	80.65	74.12	76.58	75.50
KNN	79.33	75.91	77.49	93.85	93.85	76.80	82.38
KRR (Linear)	73.09	76.06	78.06	67.95	69.20	75.77	70.35
KRR (Non-Linear)	77.31	81.99	88.37	82.76	80.51	80.74	77.82
LSTWSVM (Linear)	74.46	64.95	79.68	71.10	71.17	67.50	63.50
LSTWSVM (Non-Linear)	79.17	80.32	73.93	80.24	78.04	73.64	69.01
MPRaF-N	69.30	68.34	72.10	84.38	79.08	72.00	65.90
MPRaF-P	73.42	79.20	79.79	70.12	72.80	73.96	79.10
MPRaF-T	80.55	76.99	78.81	89.25	87.53	75.34	78.86
Neural	81.73	82.32	87.21	67.54	68.51	73.54	68.25
pinGTSVM	64.76	63.91	71.63	40.62	41.34	54.36	62.84
RaF-LDA	83.19	78.10	86.60	68.02	74.80	70.67	77.14
RaF-PCA	80.83	74.91	74.58	75.54	78.01	73.66	73.32
RaF	72.41	68.90	77.12	73.87	74.88	78.80	76.43
RELSTSVM (Linear)	70.61	69.02	82.85	70.41	72.06	69.71	66.21
RELSTSVM (Non-Linear)	74.80	78.27	82.67	78.81	79.37	74.30	69.01
RVFLAE	69.20	63.42	75.13	63.63	66.25	70.43	58.52
RVFL	77.24	81.03	75.90	78.83	72.12	69.63	62.67
SVM	81.17	76.83	73.09	70.94	65.20	67.19	62.42
TBSVM (Linear)	80.96	81.52	81.66	73.71	79.65	72.67	72.22
TBSVM (Non-Linear)	78.12	87.66	78.85	84.54	87.08	82.45	87.46
TWSVM (Linear)	76.78	73.16	80.23	73.71	75.51	70.76	67.18
TWSVM (Non-Linear)	77.45	81.79	82.46	84.54	84.40	82.21	73.98

Table S-13: Precisions for Combined Matter for 500 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	75.62	75.06	81.49	75.64	70.40	75.96	65.13
KNN	74.84	71.14	77.50	75.83	75.83	69.92	75.78
KRR (Linear)	69.43	72.86	76.24	64.59	66.25	72.92	63.96
KRR (Non-Linear)	73.21	79.03	86.25	77.51	75.37	78.31	71.71
LSTWSVM (Linear)	69.00	63.07	72.29	64.18	62.74	65.72	63.80
LSTWSVM (Non-Linear)	70.68	69.05	80.06	74.73	72.87	69.52	64.72
MPRaF-N	67.58	68.98	70.85	75.31	74.68	69.81	62.92
MPRaF-P	71.45	77.73	77.59	64.83	71.02	71.94	76.78
MPRaF-T	74.83	74.04	77.42	NaN	NaN	70.17	75.12
Neural	80.20	80.61	87.24	63.92	63.95	74.52	65.20
pinGTSVM	66.27	64.56	72.88	55.76	53.73	58.57	61.69
RaF-LDA	76.31	74.96	84.23	60.50	72.53	67.73	72.27
RaF-PCA	75.44	75.47	72.06	70.89	72.67	69.89	69.25
RaF	69.23	69.02	76.00	68.88	69.63	75.95	76.56
RELSTSVM (Linear)	72.63	71.33	76.85	63.32	63.82	70.77	63.37
RELSTSVM (Non-Linear)	73.92	73.71	78.88	77.60	77.72	75.82	64.72
RVFLAE	66.69	58.72	72.63	61.13	63.67	68.33	53.98
RVFL	74.00	76.77	74.33	73.37	65.81	69.51	58.26
SVM	77.53	72.47	71.74	63.19	63.17	63.40	59.42
TBSVM (Linear)	74.73	75.41	81.32	71.95	74.37	70.58	66.27
TBSVM (Non-Linear)	75.90	84.48	79.25	81.45	82.14	78.49	83.50
TWSVM (Linear)	71.08	67.09	80.73	71.95	73.40	70.25	64.05
TWSVM (Non-Linear)	74.80	79.37	83.84	79.78	79.59	77.57	69.64

Table S-14: F-Measures for Combined Matter for 500 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	66.34	68.04	80.07	62.25	56.58	71.36	55.72
KNN	64.40	67.92	69.98	NaN	NaN	58.93	61.68
KRR (Linear)	69.55	67.96	71.66	59.90	61.44	70.89	61.19
KRR (Non-Linear)	68.76	69.24	65.33	58.08	60.13	69.29	61.39
LSTWSVM (Linear)	70.20	63.84	73.25	66.09	64.82	66.32	64.28
LSTWSVM (Non-Linear)	72.08	70.76	81.03	75.43	73.49	70.54	65.75
MPRaF-N	70.78	70.60	71.45	53.91	55.92	65.11	59.20
MPRaF-P	68.34	72.75	73.48	56.90	63.14	67.54	62.73
MPRaF-T	68.77	71.40	72.41	NaN	NaN	64.09	63.23
Neural	73.34	74.86	85.18	57.01	60.31	70.10	63.85
pinGTSVM	68.20	66.16	75.04	63.13	61.03	59.21	62.01
RaF-LDA	69.06	68.77	78.27	47.85	59.60	66.50	58.34
RaF-PCA	69.84	72.78	69.43	61.02	63.47	66.49	65.46
RaF	68.91	68.87	73.43	59.07	60.50	72.37	66.94
RELSTSVM (Linear)	73.12	71.72	77.47	65.27	65.81	71.36	63.99
RELSTSVM (Non-Linear)	74.39	74.54	79.32	77.99	78.35	76.37	65.75
RVFLAE	57.65	52.59	67.52	58.46	60.51	63.50	NaN
RVFL	69.03	67.06	74.42	66.12	60.12	66.85	56.93
SVM	72.57	66.54	68.85	55.06	62.40	63.17	58.17
TBSVM (Linear)	70.14	67.61	77.60	65.61	64.10	69.67	62.57
TBSVM (Non-Linear)	73.79	71.56	77.72	76.87	74.78	71.58	59.38
TWSVM (Linear)	67.86	62.87	77.19	65.61	64.50	68.37	63.74
TWSVM (Non-Linear)	73.92	73.25	76.02	73.82	73.73	69.35	64.48

Table S-15: G-Means for Combined Matter for 500 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	67.34	68.97	80.54	64.04	58.81	72.32	57.22
KNN	65.44	68.72	71.19	43.47	43.47	60.29	63.71
KRR (Linear)	69.97	68.66	72.39	60.90	62.44	71.59	62.10
KRR (Non-Linear)	69.29	70.67	68.05	60.54	62.13	70.60	62.71
LSTWSVM (Linear)	7.19	7.06	7.00	7.08	7.05	7.15	7.77
LSTWSVM (Non-Linear)	0.47	0.47	0.49	0.48	0.48	0.48	0.48
MPRaF-N	71.53	71.50	72.08	56.93	58.63	65.93	60.36
MPRaF-P	68.82	73.51	74.44	57.84	64.53	68.60	64.83
MPRaF-T	69.46	72.00	73.20	NaN	NaN	64.88	65.06
Neural	74.31	75.78	85.54	58.43	61.17	71.37	64.24
pinGTSVM	69.09	66.80	75.79	64.19	61.98	60.27	62.86
RaF-LDA	69.93	69.59	79.12	49.95	61.98	67.03	60.58
RaF-PCA	70.29	73.48	69.91	62.54	64.51	67.66	66.34
RaF	69.39	69.73	74.37	60.45	61.53	73.08	68.68
RELSTSVM (Linear)	6.81	6.82	6.83	6.78	6.76	6.81	7.83
RELSTSVM (Non-Linear)	0.48	0.51	0.49	0.49	0.49	0.49	0.49
RVFLAE	58.81	53.72	68.52	59.33	61.26	64.37	55.20
RVFL	69.70	68.27	74.74	67.51	61.61	67.55	57.39
SVM	73.32	67.33	69.65	56.33	63.23	63.83	58.73
TBSVM (Linear)	70.75	68.68	78.13	67.20	65.95	70.33	63.25
TBSVM (Non-Linear)	74.55	73.29	78.14	77.28	75.42	72.71	63.02
TWSVM (Linear)	68.25	63.48	77.68	67.20	66.37	69.17	64.26
TWSVM (Non-Linear)	74.39	74.16	77.27	74.40	74.30	70.44	65.55

Table S-16: AUCs for Grey Matter for 1200 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	63.57	60.49	76.64	71.23	74.30	70.01	69.46
KNN	59.69	58.61	73.18	63.71	65.26	68.16	66.20
KRR (Linear)	60.88	58.58	76.71	71.00	70.38	67.07	71.53
KRR (Non-Linear)	67.94	68.73	76.21	74.02	73.18	70.22	68.02
LSTWSVM (Linear)	58.44	62.63	75.85	66.52	63.27	65.87	65.50
LSTWSVM (Non-Linear)	55.74	54.45	73.97	81.13	81.55	66.29	70.75
MPRaF-N	62.84	63.68	73.53	70.02	75.02	66.76	73.31
MPRaF-P	64.85	63.73	69.86	72.28	71.48	68.73	70.09
MPRaF-T	67.25	65.25	75.63	69.30	69.36	64.28	72.80
Neural	62.17	60.91	73.95	62.13	65.98	61.62	66.88
pinGTSVM	58.55	55.26	70.95	68.45	70.85	63.23	68.67
RaF-LDA	65.27	56.57	78.28	69.44	64.79	67.44	66.61
RaF-PCA	66.48	57.90	77.62	69.02	74.14	67.57	61.75
RaF	65.94	63.39	71.93	71.34	71.25	67.53	68.16
RELSTSVM (Linear)	53.90	57.92	76.41	68.37	67.04	72.60	67.72
RELSTSVM (Non-Linear)	58.24	59.73	72.48	86.20	86.45	65.88	70.61
RVFLAE	55.24	53.33	70.27	58.89	62.45	59.08	65.41
RVFL	61.64	64.68	71.14	72.08	67.09	67.04	66.40
SVM	54.87	60.42	74.90	73.29	75.38	61.09	75.05
TBSVM (Linear)	63.44	64.76	79.26	71.93	73.88	69.20	73.29
TBSVM (Non-Linear)	70.55	69.18	77.47	82.89	83.30	70.93	77.88
TWSVM (Linear)	59.68	59.19	73.84	71.93	67.43	67.45	68.39
TWSVM (Non-Linear)	66.19	66.75	77.55	82.76	83.30	71.66	77.88

Table S-17: Sensitivities for Grey Matter for 1200 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	54.57	55.64	72.07	66.58	69.92	63.90	60.97
KNN	46.15	46.00	66.23	34.75	37.84	59.15	37.56
KRR (Linear)	54.17	53.58	72.31	59.50	59.92	61.63	67.98
KRR (Non-Linear)	45.77	51.84	60.72	51.81	52.81	64.54	54.78
LSTWSVM (Linear)	60.56	53.87	74.67	74.08	76.50	70.01	73.63
LSTWSVM (Non-Linear)	73.31	75.38	82.35	80.67	83.46	75.95	81.84
MPRaF-N	62.53	54.53	73.58	58.56	65.58	59.18	67.86
MPRaF-P	54.90	56.56	66.52	67.19	62.08	65.07	60.62
MPRaF-T	55.46	59.91	67.25	55.56	53.38	52.21	63.33
Neural	59.89	61.77	71.55	62.11	65.88	63.29	61.83
pinGTSVM	57.42	56.69	70.28	79.02	82.95	72.44	80.06
RaF-LDA	54.70	49.90	72.98	64.33	63.84	60.62	56.18
RaF-PCA	66.02	54.49	74.26	67.53	76.08	63.46	56.39
RaF	56.63	51.54	71.80	65.96	61.84	63.77	63.75
RELSTSVM (Linear)	75.90	70.59	80.12	79.09	79.63	66.90	80.10
RELSTSVM (Non-Linear)	77.66	77.83	83.29	79.43	81.54	73.94	83.95
RVFLAE	54.48	48.60	63.29	57.81	59.99	57.09	56.64
RVFL	58.81	64.29	63.60	68.02	62.72	56.99	62.96
SVM	49.19	53.45	69.81	62.94	66.58	66.79	69.44
TBSVM (Linear)	51.40	61.98	76.41	71.20	63.01	67.63	64.94
TBSVM (Non-Linear)	50.36	51.79	79.26	85.48	82.80	64.97	68.61
TWSVM (Linear)	54.52	55.02	67.56	71.20	65.95	72.65	68.87
TWSVM (Non-Linear)	54.21	38.83	76.40	83.98	82.80	68.79	68.61

Table S-18: Specificity's for Grey Matter for 1200 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	72.56	65.33	81.21	75.88	78.69	76.12	77.95
KNN	73.24	71.21	80.13	92.68	92.68	77.17	94.83
KRR (Linear)	67.58	63.59	81.12	82.50	80.83	72.52	75.08
KRR (Non-Linear)	90.11	85.61	91.70	96.22	93.56	75.89	81.25
LSTWSVM (Linear)	58.87	56.62	75.17	72.00	73.86	66.39	72.68
LSTWSVM (Non-Linear)	65.85	70.20	80.73	80.21	83.99	77.89	79.83
MPRaF-N	63.15	72.82	73.48	81.48	84.46	74.34	78.75
MPRaF-P	74.81	70.89	73.20	77.37	80.89	72.39	79.57
MPRaF-T	79.03	70.60	84.02	83.03	85.33	76.34	82.27
Neural	64.46	60.05	76.34	62.15	66.08	59.94	71.94
pinGTSVM	59.67	53.84	71.62	57.88	58.74	54.02	57.27
RaF-LDA	75.83	63.24	83.58	74.55	65.74	74.25	77.03
RaF-PCA	66.95	61.31	80.98	70.50	72.20	71.69	67.12
RaF	75.25	75.25	72.05	76.71	80.67	71.30	72.58
RELSTSVM (Linear)	67.19	63.38	77.97	74.59	75.28	69.31	79.17
RELSTSVM (Non-Linear)	74.87	71.44	82.33	79.69	82.02	75.74	83.14
RVFLAE	56.01	58.05	77.25	59.98	64.91	61.08	74.19
RVFL	64.46	65.08	78.69	76.13	71.45	77.09	69.85
SVM	60.55	67.38	80.00	83.63	84.18	55.38	80.66
TBSVM (Linear)	75.49	67.54	82.12	72.66	84.74	70.77	81.63
TBSVM (Non-Linear)	90.75	86.58	75.69	80.29	83.79	76.89	87.15
TWSVM (Linear)	64.84	63.35	80.13	72.66	68.91	62.25	67.91
TWSVM (Non-Linear)	78.17	94.67	78.70	81.54	83.79	74.52	87.15

Table S-19: Precisions for Grey Matter for 1200 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	68.39	62.74	77.72	72.17	76.19	72.81	74.44
KNN	70.13	66.25	77.86	83.83	83.50	69.95	88.00
KRR (Linear)	60.49	57.71	78.17	75.13	74.43	68.59	73.82
KRR (Non-Linear)	84.74	81.06	89.31	92.50	88.33	73.06	74.64
LSTWSVM (Linear)	56.76	57.24	74.15	67.33	65.93	NaN	66.24
LSTWSVM (Non-Linear)	59.57	59.09	76.17	79.34	80.65	67.14	73.47
MPRaF-N	63.22	68.62	74.17	76.95	84.38	74.10	74.70
MPRaF-P	69.63	66.87	74.38	76.87	79.38	71.84	73.93
MPRaF-T	76.52	68.12	83.14	79.22	84.13	70.22	78.05
Neural	63.63	60.39	75.83	60.07	67.48	59.73	69.13
pinGTSVM	58.06	54.13	71.21	63.91	65.47	60.11	65.58
RaF-LDA	71.08	62.04	81.30	75.15	65.41	68.06	72.95
RaF-PCA	67.06	58.78	79.53	69.92	72.59	65.48	64.94
RaF	69.84	68.88	72.76	75.18	77.33	68.42	69.46
RELSTSVM (Linear)	58.22	59.20	75.82	69.33	69.85	68.70	69.98
RELSTSVM (Non-Linear)	62.94	63.53	75.74	81.38	82.74	67.28	74.78
RVFLAE	55.19	56.17	74.26	57.52	62.43	56.67	68.43
RVFL	61.25	65.62	74.15	75.07	68.61	71.23	68.50
SVM	54.29	61.99	76.53	79.13	79.36	59.15	78.75
TBSVM (Linear)	66.49	63.43	80.25	71.45	77.54	69.51	79.71
TBSVM (Non-Linear)	88.17	84.69	77.06	80.80	84.21	73.06	84.50
TWSVM (Linear)	61.44	59.94	78.50	71.45	68.70	65.67	67.22
TWSVM (Non-Linear)	73.55	92.50	79.67	81.49	84.21	77.97	84.50

Table S-20: F-Measures for Grey Matter for 1200 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	57.93	57.06	73.10	68.43	71.01	66.81	63.71
KNN	52.62	52.10	70.19	46.47	48.91	61.96	50.05
KRR (Linear)	55.32	53.49	73.79	64.42	64.66	62.90	67.68
KRR (Non-Linear)	58.07	61.82	69.89	65.41	65.34	66.59	61.71
LSTWSVM (Linear)	57.67	58.37	74.82	68.28	67.23	65.31	67.61
LSTWSVM (Non-Linear)	60.18	60.61	76.75	80.00	81.69	69.51	74.37
MPRaF-N	60.65	58.57	71.17	64.23	70.69	62.07	68.81
MPRaF-P	58.54	58.62	67.40	69.27	67.50	65.71	64.66
MPRaF-T	61.67	61.69	71.52	63.46	63.49	57.53	67.67
Neural	60.27	59.63	72.62	59.09	63.52	59.68	63.90
pinGTSVM	55.75	53.07	70.21	69.09	72.04	63.94	69.18
RaF-LDA	59.94	51.38	75.62	65.97	63.28	62.77	60.42
RaF-PCA	64.32	55.12	74.92	66.93	72.67	63.29	58.01
RaF	59.46	55.91	69.33	68.31	67.29	64.62	65.85
RELSTSVM (Linear)	59.36	59.92	76.50	70.39	70.50	69.81	71.65
RELSTSVM (Non-Linear)	64.66	64.53	76.56	82.15	83.48	68.97	75.81
RVFLAE	52.98	50.30	66.02	55.25	58.63	NaN	59.58
RVFL	59.24	64.13	67.25	69.60	63.64	61.61	63.13
SVM	50.86	56.61	71.76	67.44	70.09	60.93	71.42
TBSVM (Linear)	56.71	60.01	77.44	70.21	68.12	66.42	68.45
TBSVM (Non-Linear)	61.48	61.32	77.37	81.60	82.11	66.94	73.70
TWSVM (Linear)	56.00	56.19	71.48	70.21	65.76	66.83	65.19
TWSVM (Non-Linear)	61.12	53.27	76.69	81.35	82.11	68.34	73.70

Table S-21: G-Means for Grey Matter for 1200 features.

Methods	T-Test	ROC	Wilcoxon	Entropy	Bhattacharyya	MRMR	NCA
Het-RaF	59.63	58.10	73.98	68.90	72.02	67.57	65.64
KNN	55.02	53.96	71.11	52.11	54.07	63.22	55.86
KRR (Linear)	56.29	54.53	74.50	65.84	65.89	63.98	69.24
KRR (Non-Linear)	61.50	64.04	72.36	68.63	67.88	67.68	63.17
LSTWSVM (Linear)	10.42	10.49	10.41	10.56	10.49	10.56	10.55
LSTWSVM (Non-Linear)	0.54	0.55	0.53	0.55	0.56	0.55	0.54
MPRaF-N	61.72	60.03	72.50	65.95	72.75	64.08	70.02
MPRaF-P	60.30	60.12	68.86	70.61	69.08	67.01	65.94
MPRaF-T	63.74	62.81	73.28	65.38	66.02	59.31	69.15
Neural	61.01	60.35	73.15	60.08	65.05	60.58	64.68
pinGTSVM	56.71	54.19	70.48	70.25	73.10	65.08	70.95
RaF-LDA	61.38	53.38	76.37	67.77	63.95	63.55	62.40
RaF-PCA	65.39	55.86	75.89	67.82	73.48	63.87	59.25
RaF	61.24	57.93	70.78	69.42	68.41	65.34	66.22
RELSTSVM (Linear)	11.61	9.94	10.12	10.06	10.14	10.09	10.23
RELSTSVM (Non-Linear)	0.58	0.57	0.56	0.56	0.55	0.57	0.56
RVFLAE	53.88	51.31	67.38	56.43	59.88	56.21	61.01
RVFL	59.63	64.54	68.05	70.56	64.64	62.82	64.41
SVM	51.29	57.16	72.46	69.11	71.48	61.93	72.71
TBSVM (Linear)	57.80	61.31	77.88	70.77	69.18	67.48	70.33
TBSVM (Non-Linear)	64.98	64.41	77.76	82.36	82.80	67.94	75.10
TWSVM (Linear)	56.97	56.83	72.24	70.77	66.54	67.97	66.59
TWSVM (Non-Linear)	62.43	58.90	77.36	82.03	82.80	70.75	75.10

- learning and semi-supervised learning, in: Foundations and Trends® in Computer Graphics and Vision, Vol. 7, Now, 2012, pp. 81–227.
- [9] L. Zhang, P. N. Suganthan, Oblique decision tree ensemble via multi-surface proximal support vector machine, *IEEE Transactions on Cybernetics* 45 (10) (2014) 2165–2176.
 - [10] O. L. Mangasarian, E. W. Wild, Multisurface proximal support vector machine classification via generalized eigenvalues, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 28 (1) (2005) 69–74.
 - [11] R. Katuwal, P. N. Suganthan, L. Zhang, Heterogeneous oblique random forest, *Pattern Recognition* 99 (2020) 107078.
 - [12] M. Fernández-Delgado, E. Cernadas, S. Barro, D. Amorim, Do we need hundreds of classifiers to solve real world classification problems?, *Journal of Machine Learning Research* 15 (1) (2014) 3133–3181.
 - [13] R. Katuwal, P. N. Suganthan, An ensemble of kernel ridge regression for multi-class classification, *Procedia Computer Science* 108 (2017) 375–383.
 - [14] C. Saunders, A. Gammerman, V. Vovk, Ridge regression learning algorithm in dual variables, in: *Proceedings of the Fifteenth International Conference on Machine Learning*, 1998, p. 515–521.
 - [15] L. Zhang, P. N. Suganthan, A comprehensive evaluation of random vector functional link networks, *Information Sciences* 367 (2016) 1094–1105.
 - [16] Y. Zhang, J. Wu, Z. Cai, B. Du, S. Y. Philip, An unsupervised parameter learning model for RVFL neural network, *Neural Networks* 112 (2019) 85–97.
 - [17] A. Beck, M. Teboulle, A fast iterative shrinkage-thresholding algorithm for linear inverse problems, *SIAM Journal on Imaging Sciences* 2 (1) (2009) 183–202.
 - [18] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980.

- [19] M. Friedman, A comparison of alternative tests of significance for the problem of m rankings, *The Annals of Mathematical Statistics* 11 (1) (1940) 86–92.
- [20] M. Friedman, The use of ranks to avoid the assumption of normality implicit in the analysis of variance, *Journal of the American Statistical Association* 32 (200) (1937) 675–701.
- [21] R. L. Iman, J. M. Davenport, Approximations of the critical region of the Friedman statistic, *Communications in Statistics-Theory and Methods* 9 (6) (1980) 571–595.
- [22] P. B. Nemenyi, *Distribution-free multiple comparisons.*, Princeton University, 1963.