

Comparison of feature representations in MRI-based MCI-to-AD conversion prediction

Marta Gómez-Sancho^a, Jussi Tohka^{b,*}, Vanessa Gómez-Verdejo^{a,*}, for the Alzheimer's Disease Neuroimaging Initiative^c

^a*Department of Signal Processing and Communications, Universidad Carlos III de Madrid, Leganes, Spain*

^b*University of Eastern Finland, AI Virtanen Institute for Molecular Sciences, Kuopio, Finland*

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

Abstract

Alzheimer's Disease (AD) is a progressive neurological disorder in which the death of brain cells causes memory loss and cognitive decline. The identification of at-risk subjects yet showing no dementia symptoms but who will later convert to AD can be crucial for the effective treatment of AD. For this, Magnetic Resonance Imaging (MRI) is expected to play a crucial role. During recent years, several Machine Learning (ML) approaches to AD-conversion prediction have been proposed using different types of MRI features. However, few studies comparing these different feature representations exist, and the existing ones do not allow to make definite conclusions. We evaluated the performance of various types of MRI features for the conversion prediction: voxel-based features extracted based on voxel-based morphometry, hippocampus volumes, volumes of the entorhinal cortex, and a set of regional volumetric, surface area, and cortical thickness measures across the brain. Regional features consistently yielded the best performance over two classifiers (Support Vector Machines and Regularized Logistic Regression), and two datasets studied. However, the performance

*These two authors share the senior authorship.

**Corresponding author: vanessa@tsc.uc3m.es.

difference to other features was not statistically significant. There was a consistent trend of age correction improving the classification performance, but the improvement reached statistical significance only rarely.

Keywords: Alzheimer’s Disease, Magnetic Resonance Imaging, Brain, Machine Learning, Feature Representations

1. Introduction

Alzheimer’s Disease (AD) is a progressive neurological disorder in which the death of brain cells causes memory loss and cognitive decline. The progression of the neuropathology in AD starts long before clinical symptoms of the disease become apparent [1, 2, 3, 4, 5]. Also, the symptoms become progressively worse, and much effort has been placed on the early diagnosis of the AD. Related to this, Mild Cognitive Impairment (MCI), defined as a transitional phase from cognitive changes of normal aging to those typically found in dementia, is an important construct [6]. Subjects with MCI present a high risk of developing AD, but still, most people with MCI will not progress to dementia (or AD) even after 10 years of follow-up [7, 8]. Thus, identifying MCI subjects who convert to AD can be crucial for the effective treatment of AD.

Neuroimaging techniques have shown promise as tools for presymptomatic AD detection [9, 10]. Much research has been focused on T1-weighted Magnetic Resonance Imaging (MRI). It is one of the most widely studied imaging techniques [11] because it is completely non-invasive, highly available, inexpensive compared to positron emission tomography and has an excellent contrast between different soft tissues. Over the past few years, many potential MRI markers, such as the whole brain, hippocampal, and entorhinal cortex atrophy, have been shown to have diagnostic value [12]. Also, these markers have been used as the features for Machine Learning (ML) algorithms trying to predict MCI-to-AD conversion.

Indeed, there has been a surge of proposed ML algorithms for automatically predicting the future conversion from MCI to AD based on MRI (e.g.,

[13, 14, 15, 16]). This is partly driven by the free availability of large, high-quality datasets such as ADNI ¹. However, the principal focus has been in the development of new ML techniques, and their comparative evaluation has received much less attention. In particular, ML algorithms have used different types of feature sets extracted from MRI, including hippocampal volumes, volumes of the entorhinal cortex, cortical thickness measures, as well as voxel-based morphometry (VBM) features (e.g., [17, 18, 19, 20, 21] and [22] for a recent review). Despite that, systematic studies of advantages/disadvantages of various feature sets have been limited so far, and existing studies do not allow to make definite conclusions. To add to the confusion, high dimensional feature sets, such as cortical thickness or voxel-based morphometry, must be coupled with dimensionality reduction technique such as averaging the values within a brain region, Principal Component Analysis (PCA) or feature selection (see [23] for a review).

Existing comparisons between different feature representations do not provide a clear answer to the question we are interested in: "Is there a preferred representation of MRI for AD-conversion prediction?". There are multiple reasons for this. The comparisons have been geared to the AD vs. control classification problem ([24, 25, 26]), they have not included voxel-based representations [24, 27], they have utilized very short follow-up (18-months [27, 28]), they have been based on a single learning algorithm [27, 29] and/or have had highly unbalanced pMCI and sMCI classes (in [29] 149 of 165 MCI subjects converted during the 4-year follow-up that is in stark contrast to conversion rates reported in other analyses [8]). An early and important study [28], which we want to highlight, compared various feature representations including hippocampal volumes, cortical thickness, and VBM with and without regional averaging. No feature representation in this study managed to perform significantly better than chance. This somewhat disappointing result could be because 1) the methods were early ones, mostly geared to the much easier normal control vs. AD sub-

¹Information and data can be found at adni.loni.usc.edu

ject classification problem, 2) the dataset was smaller than the one currently
 55 available, and 3) the MCI non-converter was somewhat arbitrarily defined as a
 subject who did not convert in 18 months period. Moradi et al. [30] evaluated
 their method over the same dataset as [28] managing to obtain significantly
 better performance than the chance level, pointing to the reason 1) as the most
 significant cause of the improvement.

60 Since [28], we can find few studies of different feature representations pre-
 senting partially conflicting results. As an example, [21] found that the prog-
 nostic efficacy of hippocampus volumetry was better than combined regional
 volumetrics in two commercially available brain volumetric software packages
 for MCI conversion prediction. On the other hand, Gaser et al. [17] have
 65 demonstrated the superiority of their voxel-based brainAGE approach over the
 hippocampus volume biomarker and Westman et al. have emphasized the im-
 portance of having a complete set of regional features [27, 24]. Some researchers
 have opted to study feature selection, either supporting [20, 31] or opposing
 [32] data-driven feature selection. The comparisons of different automatic algo-
 70 rithms for hippocampal [33] and entorhinal cortex volumetry [34] have indicated
 that the algorithm-choice did not affect the classification accuracy. Intracranial
 volume adjustment to regional volumetry appears to only have subtle effects to
 the conversion prediction accuracy [35, 27]. Finally, it has been demonstrated
 that the neuropsychological test scores are the best predictors of conversion, but
 75 combining them with MRI information leads to improved prediction accuracy
 [30, 36].

To close this information gap, we asked what type of feature representations
 are the best for the MRI-based AD-conversion prediction. We used follow-up
 period of 3 years to define the AD-conversion, twice longer than in [27, 28].
 80 We evaluated the performance of various MRI features, including VBM-style,
 voxel-based features [17], coupled with feature preselection [30] or PCA-based
 dimensionality reduction, hippocampus volumes, volumes of the entorhinal cor-
 tex, and a complete set of regional volumetry, surface area, and cortical thickness
 measures extracted by FreeSurfer. This complements earlier studies which did

85 not include voxels -based representations [27, 21]. We additionally evaluated age removal [30, 37], which have been found to improve the prognostic efficacy of ML-based MRI biomarkers. Moreover, we used two different classifiers (Support Vector Machines, SVM, and Regularized Logistic Regression, RLR) for reducing the classifier specificity of the conclusions and trained them applying a repeated
90 10-fold cross-Validation (CV) with a sound statistical inference to compare the methods, which can be seen as an improvement of separate training and test sets in [28].

2. Material and methods

ADNI data

95 Data is collected from the the Alzheimers Disease Neuroimaging Initiative (ADNI) public database, available at adni.loni.usc.edu. The ADNI initiative was launched in 2003 as a public-private partnership, led by Principal Investigator Michael W. Weiner, MD. The primary goal of ADNI has been to test whether serial magnetic resonance imaging (MRI), positron emission tomography (PET), other biological markers, and clinical and neuropsychological
100 mography (PET), other biological markers, and clinical and neuropsychological assessment can be combined to measure the progression of mild cognitive impairment (MCI) and early Alzheimers disease (AD). For up-to-date information, see www.adni-info.org.

ADNI material considered in this work included all subjects from ADNI1
105 for whom baseline MRI data (T1-weighted MP-RAGE sequence at 1.5 Tesla, typically 256 x 256 x 170 voxels with the voxel size of approximately 1 mm x 1 mm x 1.2 mm) and sufficient follow-up information were available. We focused on the classification of MCI individuals based on their future diagnosis (AD or not AD) and, therefore, MRI scans were all obtained at the baseline visit.

110 Two flavors of this dataset were evaluated. The first dataset, Quality Control (QC) dataset, included 183 MCI subjects whose FreeSurfer 4.3 MRI segmentations had passed the complete quality control. The second one, Non QC dataset, included the complete dataset of 264 MCI subjects without any quality control.

The reason for evaluating the two different sets was to study if the quality control yielded an improvement in the data analysis. Note that the QC dataset
115 was a subset of the non-QC dataset.

Following [30], a subject was considered as a progressive MCI (pMCI) if diagnosed as MCI at baseline and the diagnosis changed to AD during the 3-year follow-up period. The subject was considered as a stable MCI (sMCI)
120 if diagnosed as MCI at baseline and the diagnosis remained as MCI during the follow-up. The minimum length of follow-up was 3 years and the subject was excluded from the study if she converted after the 3 year follow-up, the diagnosis fluctuated after the 3-year follow-up period, or less than 3-years of follow-up information was available. Table 1 lists the main characteristics of the subjects of each dataset and the list of Roster IDs of the included subjects
125 and their diagnostic categories are available in the supplement.

Table 1: Demographics of the two flavors of the dataset (QC and Non-QC) used in this work. The NC and AD subjects' data were not used directly in the learning algorithms. The NC subjects were used for the age-correction. The AD and NC subjects were used in Moradi-method for feature selection.

	QC dataset				Non-QC dataset			
	sMCI	pMCI	AD	NC	sMCI	pMCI	AD	NC
No. subjects	73	110	126	182	100	164	200	231
Males / Females	46/27	58/52	65/61	91/91	66/34	97/67	103/97	119/112
Age range	59-88	55-89	55-91	60-90	57-88	55-89	55-91	60-90

2.1. Image preprocessing

Table 2 details the feature representations we investigated and their respective number of features. Hippocampus volumes consisted of left and right hippocampal volumes. Hippocampus + Entorhinal volumes consisted of left and
130 right volumes of hippocampus and left and right volumes of entorhinal cortex. We considered both raw volumes as well as volumes normalized by the intracranial volume (ICV) as it is still unclear if the normalization by ICV is beneficial for the prediction task [35, 27]. Region based features included a complete set

135 of 257 regional cortical thickness, surface area and volume measures provided by FreeSurfer ², ³. We note that this set of features included also ICV.

Freesurfer 4.3 software was employed for the extraction of hippocampus and entorhinal cortex volumes as well as region features. Particularly, the FreeSurfer 4.3 processing results available at the ADNI website were used (UCSF Cross-sectional Freesurfer version 4.3), and the description of the pipeline and the QC procedure can be found at ⁴. The rationales for using the processing results provided by ADNI was to ensure that the processing pipeline was a standard one, the processing results are readily available to other researchers, and the quality control, independent from the authors of this study, has been performed. We 140 note that albeit different versions of FreeSurfer can result in different segmentations, the classification results based on different software versions have been found to be the same [34].

Voxel-Based Morphometry (VBM) based features consisted of 29852 gray matter density values from the VBM style preprocessing by the VBM8 software. 150 In brief, the MRIs were preprocessed into gray matter tissue images in the stereotactic space as described in [17, 30], smoothed with the 8-mm FWHM Gaussian kernel, resampled to 4 mm spatial resolution, and masked into 29852 voxels. In the Moradi set of features, VBM features were further processed through the feature selection method of [30]. This method applies MRIs of AD 155 and NC subjects to select features for MCI classification through a repeated application of the elastic net penalized linear regression. We applied the ADNI data from 231 (182) normal controls and 200 (126) AD subjects for this feature selection with the non-QC (QC dataset). We reduced the number of VBM

²<http://surfer.nmr.mgh.harvard.edu/>

³Originally, this set included 274 measures. We selected a subset of 256 regions from the aforementioned 274 measures discarding the regions that presented missed data. A more detailed description of the 256 features is provided in <https://github.com/MartaGomez/Regions-list-/wiki/Regions-list>.

⁴<https://adni.bitbucket.io/reference/docs/UCSFFRESFR/UCSFFreeSurferMethodsSummary.pdf>

features also using principal component analysis (PCA). For this, we retained
 160 the PCA components that explained 99 % of the variance.

Table 2: Summary of the sets of features considered in this study. Note that the number of Moradi and PCA voxel features was dataset dependent.

Feature Set	Number of features
Hippocampus volumes	2
Hippocampus + Entorhinal volumes	4
Region	257
Voxel	29852
Moradi	525 (non-QC); 431 (QC)
PCA Voxel	225 (non-QC); 157 (QC)

We further evaluated the representations with and without the age correction. The age correction may be important as the effects of normal aging on the brain structure partially overlap with the effects of AD [38, 37]. We applied the age correction procedure of [30]. This method estimates the age effect by a
 165 linear regression for each feature separately based on the MRIs of normal controls (231 normal controls with the age range from 55 to 90 years of ADNI) and then adjusts the features of the MCI subjects based on the estimated model.

2.2. Validation and test procedure

For the implementation and evaluation of the classification methods, we performed a repeated and nested 10-fold Cross Validation (CV). In the outer CV
 170 loop, data was split in 10 different folds from which one fold at time was designated as the test fold (for performance evaluation) and the nine remaining folds were used for classifier training. The train/test cycle was repeated with each fold once as the test fold. In the inner CV loop, each train fold was, itself,
 175 split into 10 validation folds from which one part was used to select the classifier hyperparameters. The optimal hyperparameters were selected evaluating either the classification accuracy (ACC, number of correctly classified samples

over the total number of samples) or the Area Under the receiving operating Curve (AUC) [39]. The nested CV was repeated 10 times, each with a different
 180 randomly selected folding scheme, to minimize the effect of a particular folding scheme to the results. Also, the hypothesis test we used to compare different representations requires the repeated use of CV.

To study the classifier performance, we considered several metrics: AUC, ACC, Sensitivity (SEN, number of correctly classified pMCI subjects divided by
 185 the total number of pMCI subjects and Specificity (SPE, number of correctly classified sMCI subjects divided by the total number of sMCI subjects). We selected AUC as our principal performance measure as it is insensitive to the class-imbalance whereas ACC can be strongly affected by the class-imbalance.

2.3. Classifiers

To evaluate each feature set, we considered two types of widely used supervised learning classifiers: Support Vector Machine (SVM) [40] and elastic-net Regularized Logistic Regression (RLR) [41]. Accessible description of these learning methods can be found in [42]. For the SVM implementation, we used the Python open source library Scikit-learn (<http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>), which is based on the
 195 LIBSVM implementation (<https://www.csie.ntu.edu.tw/~cjlin/libsvm/>). For the RLR classifiers, we applied the GLMNET Python library (https://web.stanford.edu/~hastie/glmnet_python/), which solves the resulting penalized optimization problem by a coordinate descent algorithm. We note
 200 that both of these learning algorithms tolerate high-dimensional data via regularization and are therefore suited for the cases where the number of features is higher than the number of subjects. Especially, elastic-net includes an L1-penalty, which leads to feature selection embedded to the classifier learning [43]. A large majority of supervised learning techniques have utilized these learning
 205 algorithms [22] and a comparison of different classification algorithms MCI-to-AD prediction is available in [44]. We note that we did not include Random Forests [45] as these are not straight-forwardly suitable for high-dimensional

small-sample problems and the computation time and memory requirements for nearly all implementations would be prohibiting for the voxel-based features (however, see [46]).

For the case of the SVM classifier, we decided to use the linear SVM (we have also analyzed the possibility of using a RBF (Radial Basis Function) kernel, however, experimental results showed similar performances). In this way, we had to select only the soft margin parameter, C , whose value was explored among the set $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2, 10^3\}$ (see [40] for notation). Despite considering the linear SVM, its implementation was carried out in the dual space, precomputing a linear kernel; in this way, we simplified the calculations and reduced the computation time with the high dimensional feature representations, such as VBM ones.

For the RLR classifier, using the notation of [41], we set the parameter of the elastic net α to 0.5, just in between lasso ($\alpha=1$) and ridge ($\alpha=0$) regularization. The principal regularization parameter of the RLR (λ), which sets the balance between the regularization and the data terms, was chosen among the set of values $\{10^{-10}, 10^{-4}, 10^{-3}, 5 \cdot 10^{-3}, 10^{-2}, 5 \cdot 10^{-2}, 10^{-1}, 5 \cdot 10^{-1}\}$.

Finally, as a step prior to training the classifiers, we normalized the data by removing its mean and scaling it to the unit variance.

2.4. Statistical test

To compare the AUC values provided by different approaches, we applied the corrected resampled t-test [47]. The problem in applying standard statistical methodology, such as uncorrected t-test to assess the differences between AUCs is that $r \times k$ AUC values from in a k -fold CV repeated r times are not statistically independent. Instead, the corrected resampled t-test assumes dependency among the AUCs in a k -fold CV repeated r times and, therefore, it allows to statistically compare two mean AUC values by correcting the variance estimation. The corrected resampled t-test can be seen as an improvement over the 5x2 CV of [48] and McNemar's test for the classification accuracy [47]. Although the test was developed for the classification accuracy, it is as well

applicable for testing the differences between AUCs.

To describe the test formally, let n_1 and n_2 , respectively, denote the number
of instances used for training and testing in each fold, and a_{ij} and b_{ij} represent
the AUCs of the i -th fold and j -th run of the method A and B with $i =$
 $1, \dots, k$ and $j = 1, \dots, r$. Denoting the estimated mean and variance values of the
differences between methods A and B by \hat{m} and $\hat{\sigma}^2$ i.e.,

$$\hat{m} = \frac{1}{kr} \sum_{i=1}^k \sum_{j=1}^r a_{ij} - b_{ij} \quad (1)$$

$$\hat{\sigma}^2 = \frac{1}{kr-1} \sum_{i=1}^k \sum_{j=1}^r (a_{ij} - b_{ij} - \hat{m})^2 \quad (2)$$

we can estimate the statistic of the test, t , as:

$$t = \frac{\hat{m}}{\sqrt{\left(\frac{1}{kr} + \frac{n_2}{n_1}\right) \hat{\sigma}^2}} \quad (3)$$

The statistic t follows a student's t-distribution with $kr - 1$ degrees of freedom.

In our case, $r = k = 10$.

3. Results

Tables 3, 4 and 5 show, respectively, the results for the QC dataset and
non-QC datasets using the AUC for model selection and the results of the non-
QC dataset when the best model was selected using ACC. In particular, each
table includes for the SVM and RLR classifiers the values of AUC (area un-
der the curve), ACC (accuracy), SEN (sensitivity), SPE (specificity), as well as
three p-values from hypothesis tests comparing the AUCs: p_{Age} (comparing age
removed features vs. non age removed features), p_{Hippo} (comparing hippo-
ampus features with the remaining features for the age removed case) and p_{Class}
(comparing SVM vs. LR results over the same set of features).

The AUC values of region features were the highest in all the experiments.
However, the performance improvement over the hippocampus feature set, which
was our baseline, did not reach the statistical significance and these improved

Table 3: Cross-validated performance measures with the QC dataset using AUC as the model selection criterion. Hippocampus (Hippo. and entor. vol.) volumes ICV refers to hippocampus (hippocampus + entorhinal) volumes normalized by the ICV.

Classifier	Feature Set	Age Removal	AUC	ACC	SEN	SPE	P_{Age}	P_{Hippo}	P_{Class}
SVM	Hippocampus	No	73.50 %	63.15 %	93.00 %	18.16 %			0.873
SVM	Hippocampus	Yes	77.31 %	66.42 %	90.54 %	30.20 %	0.052		0.819
SVM	Hippocampus vol. ICV	No	69.61 %	67.21 %	81.36 %	45.82 %			0.690
SVM	Hippocampus vol. ICV	Yes	72.67 %	68.31 %	82.36 %	47.11 %	0.124	0.169	0.416
SVM	Hippo. and entor.	No	75.91 %	63.69 %	94.55 %	17.14 %			0.314
SVM	Hippo. and entor.	Yes	78.59 %	61.77 %	97.82 %	7.29 %	0.055	0.492	0.285
SVM	Hippo. & entor. vol. ICV	No	72.66 %	68.26 %	83.82 %	44.79 %			0.850
SVM	Hippo. & entor. vol. ICV	Yes	74.50 %	68.86 %	81.64 %	49.66 %	0.325	0.399	0.505
SVM	Voxel features	No	63.19 %	60.51 %	91.36 %	13.98 %			0.649
SVM	Voxel features	Yes	66.67 %	61.65 %	91.82 %	16.20 %	0.175	0.042	0.923
SVM	PCA VF	No	63.50 %	59.88 %	92.73 %	10.23 %			0.604
SVM	PCA VF	Yes	65.53 %	60.75 %	92.64 %	12.54 %	0.535	0.035	0.509
SVM	Moradi features	No	71.92 %	63.54 %	92.63 %	19.66 %			0.025
SVM	Moradi features	Yes	75.08 %	62.32 %	97.10 %	9.86 %	0.243	0.610	0.001
SVM	Region features	No	74.06 %	64.67 %	92.27 %	22.91 %			0.739
SVM	Region features	Yes	77.34 %	69.29 %	88.45 %	37.27 %	0.241	0.990	0.759
LR	Hippocampus	No	73.38 %	68.41 %	82.00 %	47.86 %			
LR	Hippocampus	Yes	77.19 %	71.31 %	83.91 %	52.43 %	0.042		
LR	Hippocampus vol. ICV	No	69.80 %	66.57 %	75.16 %	38.88 %			
LR	Hippocampus vol. ICV	Yes	72.25 %	68.32 %	84.64 %	43.68 %	0.253	0.138	
LR	Hippo. and entor.	No	72.52 %	66.85 %	80.82 %	45.73 %			
LR	Hippo. and entor.	Yes	74.13 %	68.05 %	81.64 %	47.55 %	0.085	0.767	
LR	Hippo. & entor. vol. ICV	No	71.72 %	68.59 %	84.47 %	42.60 %			
LR	Hippo. & entor. vol. ICV	Yes	73.60 %	68.90 %	83.51 %	45.00 %	0.332	0.348	
LR	Voxel features	No	64.40 %	61.91 %	76.18 %	40.43 %			
LR	Voxel features	Yes	66.94 %	63.57 %	77.36 %	40.43 %	0.268	0.029	
LR	PCA VF	No	61.95 %	59.12 %	78.73 %	26.68 %			
LR	PCA VF	Yes	63.40 %	61.33 %	80.82 %	31.91 %	0.600	0.006	
LR	Moradi features	No	65.83 %	63.28 %	75.00 %	42.73 %			
LR	Moradi features	Yes	67.89 %	65.52 %	77.54 %	57.50 %	0.348	0.038	
LR	Region features	No	74.57 %	70.67 %	81.91 %	53.70 %			
LR	Region features	Yes	77.91 %	71.93 %	81.45 %	57.50 %	0.171	0.839	

Table 4: Cross-validated performance measures with the non QC dataset using AUC as the model selection criterion. Hippocampus (Hippo. and entor. vol.) volumes ICV refers to hippocampus (hippocampus + entorhinal) volumes normalized by the ICV.

Classifier	Feature Set	Age Removal	AUC	ACC	SEN	SPE	<i>PAge</i>	<i>PHippo</i>	<i>PClass</i>
SVM	Hippocampus volumes	No	70.29 %	63.09 %	96.02 %	9.10 %			0.708
SVM	Hippocampus volumes	Yes	75.57 %	63.94 %	95.39 %	12.30 %	0.047		0.398
SVM	Hippocampus vol. ICV	No	69.56 %	68.25 %	85.60 %	39.80 %			0.987
SVM	Hippocampus vol. ICV	Yes	72.28 %	69.46 %	85.3%	43.50 %	0.172	0.181	0.807
SVM	Hippo. and entor. vol.	No	73.23 %	64.57 %	98.16 %	3.50 %			0.358
SVM	Hippo. and entor. vol.	Yes	76.05 %	63.73 %	96.86 %	9.50 %	0.088	0.744	0.253
SVM	Hippo. & entor. vol. ICV	No	71.75 %	68.55 %	84.36 %	42.60 %			0.948
SVM	Hippo. & entor. vol. ICV	Yes	73.77 %	69.60 %	82.95 %	47.90 %	0.194	0.495	0.756
SVM	Voxel based features	No	68.55%	62.11 %	97.10 %	4.80 %			0.695
SVM	Voxel based features	Yes	69.55%	63.79 %	96.55 %	10.10 %	0.659	0.169	0.489
SVM	PCA VF	No	68.41 %	62.30 %	96.19 %	6.80 %			0.328
SVM	PCA VF	Yes	69.13 %	63.90 %	95.12 %	12.70 %	0.754	0.149	0.549
SVM	Moradi features	No	73.26%	64.20%	96.04 %	12.00 %			0.078
SVM	Moradi features	Yes	75.72%	63.40 %	96.89 %	8.50 %	0.180	0.965	0.288
SVM	Region features	No	73.11%	65.29 %	90.07 %	24.60 %			0.123
SVM	Region features	Yes	76.89%	66.94 %	95.30 %	20.50 %	0.091	0.680	0.101
LR	Hippocampus volumes	No	70.95 %	67.00 %	88.04 %	32.40 %			
LR	Hippocampus volumes	Yes	74.95 %	67.68 %	86.34 %	37.00 %	0.046		
LR	Hippocampus vol. ICV	No	69.60 %	67.94 %	87.91 %	35.20 %			
LR	Hippocampus vol. ICV	Yes	72.37 %	69.44 %	88.76 %	37.80 %	0.148	0.272	
LR	Hippo. and entor. vol.	No	72.59 %	69.72 %	85.88 %	43.20 %			
LR	Hippo. and entor. vol.	Yes	75.31 %	70.61 %	84.98%	47.00 %	0.094	0.801	
LR	Hippo. & entor. vol. ICV	No	71.72 %	68.59%	84.47 %	42.60 %			
LR	Hippo. & entor. vol. ICV	Yes	73.60 %	68.90 %	83.51 %	45.00 %	0.210	0.603	
LR	Voxel based features	No	69.46 %	66.63 %	83.42 %	39.10 %			
LR	Voxel based features	Yes	71.34 %	66.99 %	82.68 %	41.30 %	0.370	0.394	
LR	PCA VF	No	66.02 %	64.50 %	86.28 %	28.80 %			
LR	PCA VF	Yes	67.58 %	64.68 %	87.02 %	28.10 %	0.536	0.058	
LR	Moradi features	No	69.94%	67.77 %	83.81 %	41.50 %			
LR	Moradi features	Yes	74.04 %	70.84 %	86.79 %	44.70 %	0.068	0.798	
LR	Region features	No	76.38 %	71.27 %	85.97 %	47.10 %			
LR	Region features	Yes	79.58 %	71.73 %	84.07 %	51.50 %	0.120	0.060	

Table 5: Cross-validated performance measures with the non-QC dataset using ACC as the model selection criterion

Classifier	Feature Set	Age Removal	AUC	ACC	SEN	SPE	<i>P_{age}</i>	<i>P_{Hippo}</i>	<i>P_{Class}</i>
SVM	Hippocampus volumes	No	70.51 %	67.35 %	85.06 %	38.30 %			0.330
SVM	Hippocampus volumes	Yes	74.99 %	68.86 %	81.91 %	47.40 %	0.026		0.759
SVM	Hippo. and entor. vol.	No	72.43 %	69.03 %	81.36 %	48.8 %			0.968
SVM	Hippo. and entor. vol.	Yes	75.40 %	71.71 %	82.20 %	53.50 %	0.065	0.777	0.880
SVM	Voxel based features	No	67.10%	61.71 %	79.86 %	32.10 %			0.939
SVM	Voxel based features	Yes	68.35%	62.46 %	80.59 %	32.90 %	0.506	0.106	0.548
SVM	PCA VF	No	66.78 %	62.15 %	79.56 %	33.70 %			0.621
SVM	PCA VF	Yes	67.98 %	63.14 %	80.18 %	35.30 %	0.613	0.085	0.784
SVM	Moradi features	No	72.85%	68.93%	84.49 %	43.40 %			0.356
SVM	Moradi features	Yes	75.00%	70.09 %	83.99 %	47.30 %	0.292	0.997	0.650
SVM	Region features	No	72.55 %	69.16%	82.73 %	46.90 %			0.122
SVM	Region features	Yes	75.98%	71.01 %	86.94 %	44.90 %	0.105	0.763	0.076
LR	Hippocampus volumes	No	70.96 %	66.28 %	88.92 %	29.10 %			
LR	Hippocampus volumes	Yes	74.85 %	69.12 %	84.10 %	44.50 %	0.050		
LR	Hippo. and entor. vol.	No	72.40 %	69.55 %	85.50 %	43.30 %			
LR	Hippo. and entor. vol.	Yes	75.50 %	70.50 %	84.68 %	47.20 %	0.056	0.650	
LR	Voxel based features	No	67.33 %	64.81 %	80.76 %	38.60 %			
LR	Voxel based features	Yes	69.95 %	66.15 %	80.21 %	43.10 %	0.235	0.244	
LR	PCA VF	No	65.42 %	63.52 %	87.46 %	24.30 %			
LR	PCA VF	Yes	67.35 %	64.96 %	87.76 %	27.60 %	0.513	0.059	
LR	Moradi features	No	71.08 %	68.60 %	85.35 %	41.10 %			
LR	Moradi features	Yes	74.10 %	70.42 %	85.94 %	45.00 %	0.125	0.836	
LR	Region features	No	75.91 %	71.01 %	84.52 %	48.80 %			
LR	Region features	Yes	79.41 %	72.07 %	84.24 %	52.10 %	0.123	0.717	

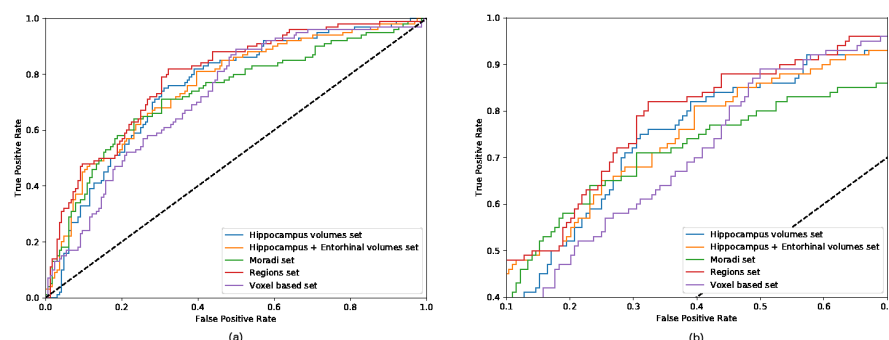


Figure 1: ROC curves corresponding to distinct the features sets used in RLR classification with the non-QC dataset. The age effect was removed.

AUCs need to be interpreted with care. In the particular case of the non-QC dataset and the RLR classifier, the regions feature set produced significantly higher AUC than hippocampus volumes.

Figure 1 depicts the ROC curves for the different feature sets under study for the RLR classifier in the non-QC dataset. Focusing on the center of these curves (see the panel 1b), we can corroborate that the region feature set appeared superior, but the performance differences were small. To avoid crowding, the ROCs of the PCA voxel feature set was not visualized as it always performed worse than the voxel features without PCA. For similar reason, the figure only displays the ROC curves of raw Hippocampus and Hippocampus + Entorhinal cortex volumes and not the ICV-normalized ones. The same principle will be followed in later figures.

Regarding the use of two different classifiers, differences between AUCs of SVM and RLR were not significant. However, SVM yielded low specificity values and the relation between SPE and SEN was more balanced with the LR classifier. Because of this we studied whether the use of AUC as the model selection criteria contributed to this imbalance with the SVM classifier. Using ACC as a model selection criterion avoids this SPE/SEN imbalance, as can be seen in Figure 2 where the specificity values are compared between ACC and AUC based model selection. As the comparison of Tables 4 and 5 reveals, the

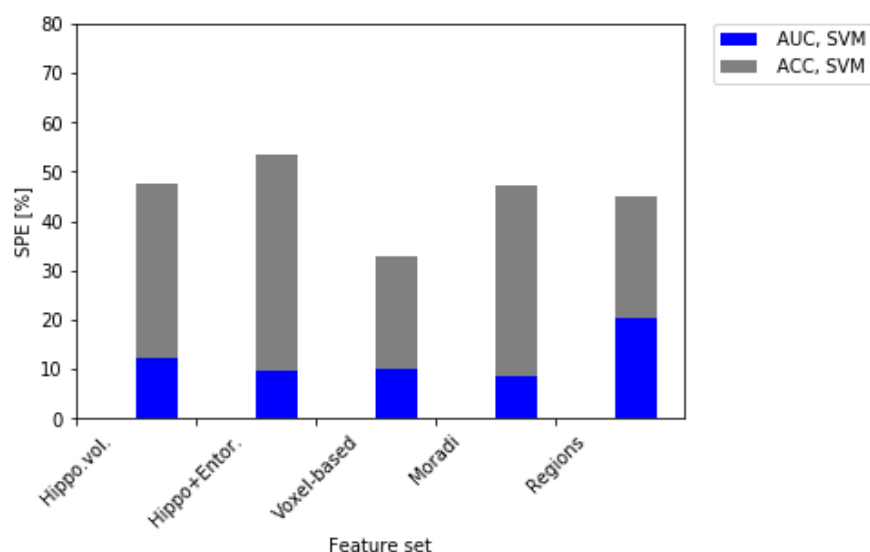


Figure 2: Specificity values of SVM classifiers when AUC and ACC were used for model selection. The models selected with ACC resulted in specificity values close 50 % whereas the models selected with AUC resulted in very low specificity values.

final AUC values did not markedly differ between the two model selectors.

We evaluated the effects of age removal on the feature sets. For this purpose, Figure 3 shows a detailed analysis of the advantages of removing the age effects. As a result, classification scores improved for every age removed effects feature set (see the panel 3c). However, as visible in Tables 3 and 4, significant improvement (p-value < 0.1) was observed only for hippocampus and hippocampus + entorhinal volume feature sets.

The differences between the AUCs of raw and ICV-normalized hippocampus and hippocampus + entorhinal volumes were not significant. Surprisingly, the raw volumes performed slightly better in terms of AUC within each dataset. However, this result agrees with findings in [35, 27] and it is not central for the purposes of this work to analyze the potential reasons for this result.

Finally, Figure 4 shows the differences between QC and non QC datasets when age effects were removed. As expected Hippocampus and Hippocam-

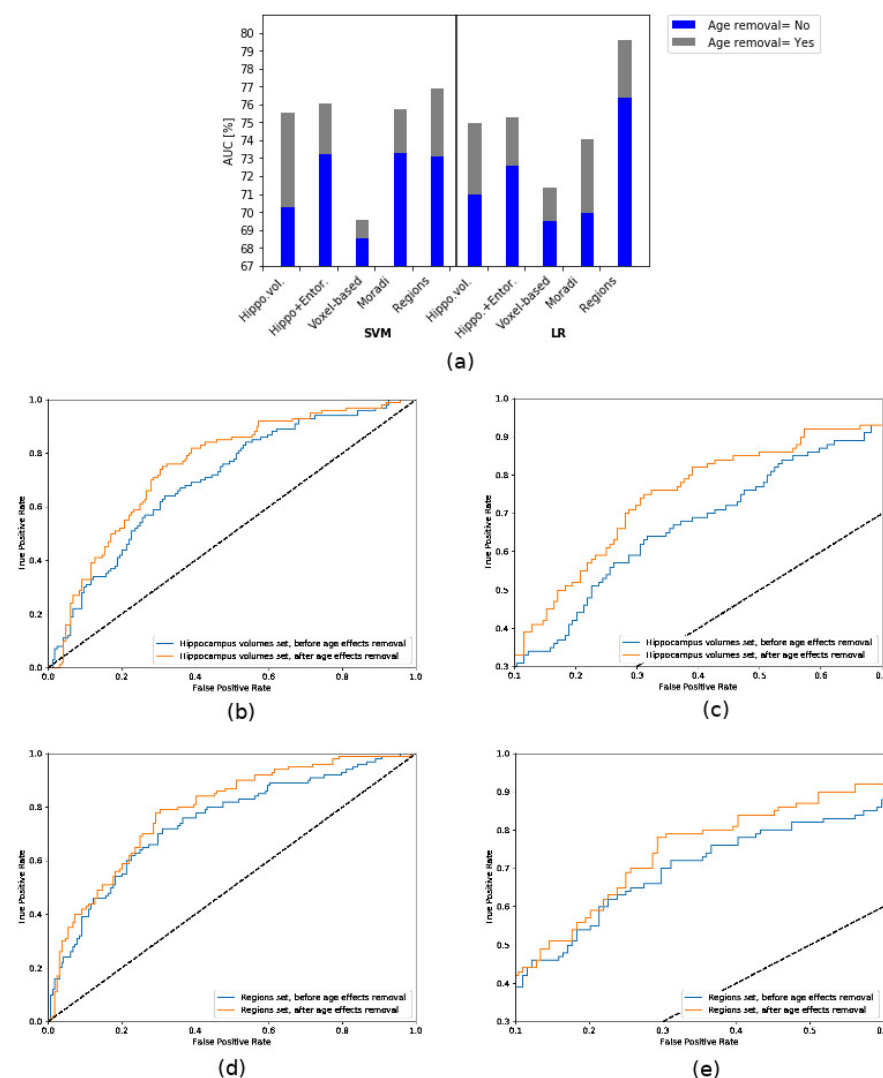


Figure 3: Analysis of age removal effects: (a) AUC comparison for different feature sets and both classifiers; (b) and (c) ROC curves for RLR classifier using hippocampus volumes; (d) and (e) ROC curves for LR classifier using region features. Age removal improved predictions in all cases.

pus plus Entorhinal volumes were benefited from the Quality Control process, whereas remaining features sets resulted in better performances when all the available data were used.

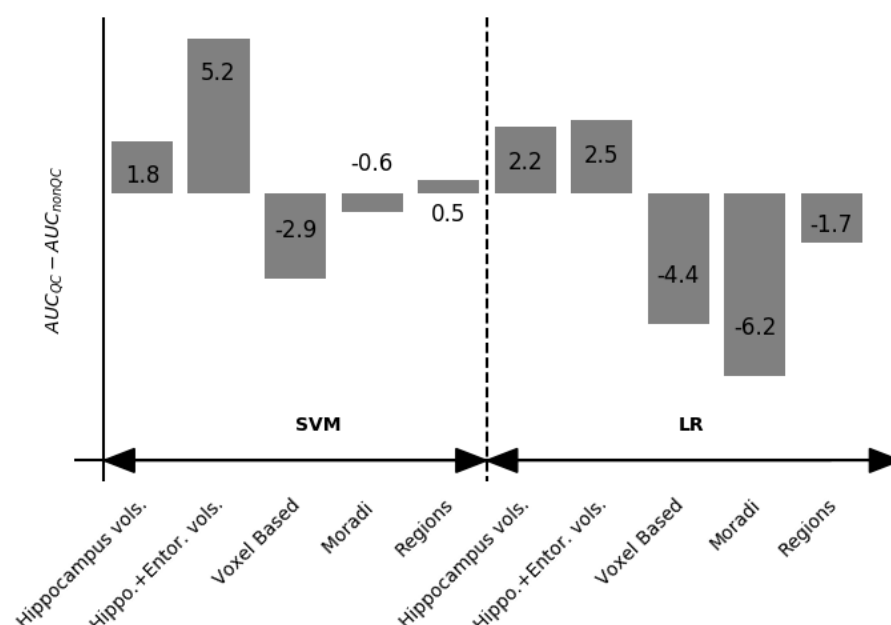


Figure 4: Differences between the AUC values with the QC dataset and the non-QC dataset for SVM (left) and RLR (right).

4. Discussion

In this work, we compared six different feature representations of MRI for predicting the AD conversion in MCI subjects. The feature sets we studied varied from high dimensional feature sets produced by VBM via regional cortical thickness, surface area, and volumetry to simple and easily interpretable features such as hippocampus and entorhinal cortex volumes (see Table 2). We addressed the feature representations using two learning algorithms, SVM and RLR, and with several metrics, AUC, ACC, SEN and SPE, that gave a reliable insight into the relative performance of different feature sets. AUC was selected as the principal figure of merit, due to its insensitivity to the class imbalance (note that the datasets contained twice the number of pMCIs (subjects who converted to AD) compared to sMCIs (subjects who remained as MCIs)). The evaluation process was carried out with a nested fold CV repeated 10 times

ensuring the insensitivity of the conclusions to random train/test division of the holdout method used previously [28]. Selecting the parameters of the classifiers inside nested CV ensures that there are no biases towards particular feature representations due to arbitrarily selected classifier parameters.

315 We found that age-corrected *regions feature set* (see <https://github.com/MartaGomez/Regions-list-/wiki/Regions-list> for a detailed description) outperformed the remaining feature sets, specifically in AUC, even though the improvement did not reach statistical significance. This result suggests that regions based features were equal or better predictors than the left and right
320 hippocampal volumes (HV) alone (which were included in the region feature set). This is interesting as a recent study [21] concluded that HV had the highest AUC among a set of individual regional volume features and was better in terms of the prognostic efficacy of combining various volumetrics. Their experimental setting was similar to the one analyzed here, however, with three
325 main differences. First, removing age related effects from MRI data was not considered; second, the set of pMCI patients was about half of ours; and, third, the combined volumetric analysis did not consider measures such as surface area or cortical thickness. This can explain the improvement in the best classification accuracy from 69 % of [21] to 80 % in the present study.

330 Voxel-based representations did not perform well in this study when coupled with standard feature reduction techniques (elastic-net or PCA). This was in contrast to a recent data-analysis competition, where the goal was to classify subjects into NC, MCI, and AD categories based on MRI [26]. However, as multiple factors have effect to a performance of an approach in a data analysis
335 competition, definite conclusions on feature representations cannot be made based on such competitions. However, also in our own experience, voxel-based methods, coupled with elastic-net feature selection, perform well in classifying between NC and AD or NC and MCI [31]. These discrepancies may suggest that NC vs. MCI (or AD) classification and AD-conversion prediction have different
340 characteristics. Further, we note that feature pre-selection based on AD and NC data suggested by Moradi et al. [30] improved the conversion prediction

accuracy markedly.

Retico et al. found that the voxel based VBM features best discriminate between sMCI and pMCI after applying Recursive Feature Elimination (RFE) [20]. However, again, the maximum accuracy in [20] was much lower than the accuracies in the present study and pMCI vs. sMCI classifiers were trained only using AD and NC subjects that may explain this. Additionally, the statistical framework was incomplete as no hypothesis testing was done and the exact definition of stable MCI class remained unclear. Other works, such as [18], concluded that the combination of different feature representations resulted into a better classification accuracy than one representation alone. Again, the classification accuracies were lower than in the present work. Moreover, [18] selected classifier hyperparameters based on test data that may cause upward bias in the reported accuracies [15].

It is important to point out that while our classification accuracies were better than those in the studies reviewed above, the performance measures are not directly comparable because different definitions of pMCI and sMCI. In fact, this is a problem that complicates the comparison of ML methods for this particular application and it is reviewed in further length in [22]. Namely, the definition of sMCI subject based on a certain cutoff (say 3 years) is problematic as this simple criterion would place a subject who received an AD diagnosis 4 years after the baseline visit into the sMCI category. Our view is that this would create unrealistic heterogeneity into the sMCI class and therefore tracking subjects' status after the cutoff is necessary (if possible). We have populated our sMCI category based on all the information available by ADNI.

Regarding the used ML methods, RLR provided, in general, similar AUC values than SVM, but had an advantage of higher specificity (it classified sMCI cases much better than the SVM did). SVM had a tendency of overpopulating the pMCI class. However, in the case of SVM, low specificity seemed to depend on the using AUC as the criterion for the hyperparameter selection. The values in Table 5 reveal how selecting the hyperparameters instead through ACC resulted in an overall improvement of specificity with a small loss of sensitivity.

This is an interesting phenomenon, as it seems to be a problem of a specific class of learning algorithms, which invites further research. However, as this
 375 issue is not central to the goals of this work, we do not analyze it further.

There were no significant differences between the classification accuracies or AUCs obtained with non-QC and QC datasets. However, the small differences between the two datasets were as expected as shown in Figure 4. For Hippocampus and Hippocampus and Entorhinal volumes, the QC was moderately useful
 380 whereas for the Moradi and Voxel based features it was moderately detrimental. This is as expected since the QC was based on Freesurfer segmentations (as Hippocampus and Entorhinal volumes) but the voxel-based and Moradi features were not. Interestingly, for region based features (also based on Freesurfer segmentation), the QC seemed not to influence the performance of the classifier.

385 It is remarkable that the age removal seem to be a key for better performances. As Figure 3 illustrates, age removal always led to better classification performances, although the improvements were not always statistically significant. This agrees with a recent work of [31] which demonstrated the same for NC vs. MCI classification.

390 5. Conclusion

This paper evaluated the performance of various types of MRI features for the future AD conversion prediction and it also analyzed the performance of each feature set over two classifiers (Support Vector Machines and Regularized Logistic Regression) and with and without applying an age correction process.

395 Experimental results showed that regional features consistently yielded the best performance, although the performance difference to other features was not statistically significant. Besides, the age removal seemed to be a key for better performances, but the improvement reached statistical significance only rarely.

Acknowledgments

400 Data collection and sharing for this project was funded by the Alzheimer's
Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant
U01 AG024904) and DOD ADNI (Department of Defense award number W81XWH-
12-2-0012). ADNI is funded by the National Institute on Aging, the National In-
stitute of Biomedical Imaging and Bioengineering, and through generous contri-
405 butions from the following: AbbVie, Alzheimers Association; Alzheimers Drug
Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-Myers
Squibb Company; CereSpir, Inc.; Eisai Inc.; Elan Pharmaceuticals, Inc.; Eli
Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affili-
ated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen
410 Alzheimer Immunotherapy Research & Development, LLC.; Johnson & Johnson
Pharmaceutical Research & Development LLC.; Lumosity; Lundbeck; Merck &
Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Neurotrack Tech-
nologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging;
Servier; Takeda Pharmaceutical Company; and Transition Therapeutics.

415 J. Tohka's work was supported by the Academy of Finland and V. Gómez-
Verdejo's work has been partly funded by the Spanish MINECO grant TEC2014-
52289R and TEC2016-81900-REDT/AEI.

References

- [1] H. Braak, E. Braak, Development of Alzheimer-related neurofibrillary
420 changes in the neocortex inversely recapitulates cortical myelogenesis, *Acta
neuropathologica* 92 (2) (1996) 197–201.
- [2] A. Delacourte, J. David, N. Sergeant, L. Buee, A. Wattez, P. Vermersch,
F. Ghzali, C. Fallet-Bianco, F. Pasquier, F. Lebert, et al., The biochemical
pathway of neurofibrillary degeneration in aging and Alzheimers Disease,
425 *Neurology* 52 (6) (1999) 1158–1158.

- [3] J. Morris, M. Storandt, D. McKeel, E. Rubin, J. Price, E. Grant, L. Berg, Cerebral amyloid deposition and diffuse plaques in “normal” aging evidence for presymptomatic and very mild Alzheimer’s Disease, *Neurology* 46 (3) (1996) 707–719.
- 430 [4] A. Serrano-Pozo, M. P. Frosch, E. Masliah, B. T. Hyman, Neuropathological alterations in Alzheimer Disease, *Cold Spring Harbor perspectives in medicine* 1 (1) (2011) a006189.
- [5] L. Mosconi, M. Brys, L. Glodzik-Sobanska, S. De Santi, H. Rusinek, M. J. De Leon, Early detection of Alzheimers Disease using neuroimaging, *Experimental gerontology* 42 (1) (2007) 129–138.
- 435 [6] R. C. Petersen, B. Caracciolo, C. Brayne, S. Gauthier, V. Jelic, L. Fratiglioni, Mild cognitive impairment: A concept in evolution, *Journal of internal medicine* 275 (3) (2014) 214–228.
- [7] S. J. Vos, F. Verhey, L. Frölich, J. Kornhuber, J. Wiltfang, W. Maier, O. Peters, E. Rüther, F. Nobili, S. Morbelli, et al., Prevalence and prognosis of Alzheimers Disease at the mild cognitive impairment stage, *Brain* 138 (5) (2015) 1327–1338.
- 440 [8] A. J. Mitchell, M. Shiri-Feshki, Rate of progression of mild cognitive impairment to dementia—meta-analysis of 41 robust inception cohort studies, *Acta Psychiatrica Scandinavica* 119 (4) (2009) 252–265.
- 445 [9] S. Teipel, A. Drzezga, M. J. Grothe, H. Barthel, G. Chételat, N. Schuff, P. Skudlarski, E. Cavedo, G. B. Frisoni, W. Hoffmann, et al., Multimodal imaging in alzheimer’s disease: validity and usefulness for early detection, *The Lancet Neurology* 14 (10) (2015) 1037–1053.
- 450 [10] M. W. Weiner, D. P. Veitch, P. S. Aisen, L. A. Beckett, N. J. Cairns, R. C. Green, D. Harvey, C. R. Jack, W. Jagust, J. C. Morris, et al., Recent publications from the alzheimer’s disease neuroimaging initiative: Reviewing progress toward improved ad clinical trials, *Alzheimer’s & Dementia*.

- [11] K. A. Johnson, N. C. Fox, R. A. Sperling, W. E. Klunk, Brain imaging
455 in Alzheimer Disease, Cold Spring Harbor perspectives in medicine 2 (4)
(2012) a006213.
- [12] G. B. Frisoni, N. C. Fox, C. R. Jack, P. Scheltens, P. M. Thompson, The
clinical use of structural MRI in Alzheimer Disease, Nature Reviews Neu-
rology 6 (2) (2010) 67–77.
- [13] S. Klöppel, C. M. Stonnington, C. Chu, B. Draganski, R. I. Scahill, J. D.
460 Rohrer, N. C. Fox, C. R. Jack Jr, J. Ashburner, R. S. Frackowiak, Au-
tomatic classification of MR scans in Alzheimer’s Disease, Brain 131 (3)
(2008) 681–689.
- [14] S. Adaszewski, J. Dukart, F. Kherif, R. Frackowiak, B. Draganski, A. D. N.
465 Initiative, et al., How early can we predict Alzheimer’s Disease using com-
putational anatomy?, Neurobiology of aging 34 (12) (2013) 2815–2826.
- [15] S. F. Eskildsen, P. Coupé, D. García-Lorenzo, V. Fonov, J. C. Pruessner,
D. L. Collins, Prediction of Alzheimer’s disease in subjects with mild cog-
nitive impairment from the adni cohort using patterns of cortical thinning,
470 NeuroImage 65 (2013) 511–521.
- [16] R. Casanova, C. T. Whitlow, B. Wagner, J. Williamson, S. A. Shumaker,
J. A. Maldjian, M. A. Espeland, High dimensional classification of struc-
tural MRI Alzheimers Disease data based on large scale regularization,
Frontiers in neuroinformatics 5.
- [17] C. Gaser, K. Franke, S. Klöppel, N. Koutsouleris, H. Sauer, A. D. N.
475 Initiative, et al., BrainAGE in mild cognitive impaired patients: predicting
the conversion to Alzheimers Disease, PloS ONE 8 (6) (2013) e67346.
- [18] R. Wolz, V. Julkunen, J. Koikkalainen, E. Niskanen, D. P. Zhang, D. Rueck-
ert, H. Soininen, J. Lötjönen, A. D. N. Initiative, et al., Multi-method anal-
480 ysis of MRI images in early diagnostics of Alzheimer’s Disease, PloS one
6 (10) (2011) e25446.

- [19] E. Westman, J.-S. Muehlboeck, A. Simmons, Combining MRI and CSF measures for classification of Alzheimer's Disease and prediction of Mild Cognitive Impairment conversion, *Neuroimage* 62 (1) (2012) 229–238.
- 485 [20] A. Retico, P. Bosco, P. Cerello, E. Fiorina, A. Chincarini, M. E. Fantacci, Predictive models based on Support Vector Machines: Whole-brain versus regional analysis of structural MRI in the Alzheimer's Disease, *Journal of Neuroimaging* 25 (4) (2015) 552–563.
- 490 [21] T. Tanpitukpongse, M. Mazurowski, J. Ikhen, J. Petrella, Predictive utility of marketed volumetric software tools in subjects at risk for Alzheimer Disease: Do regions outside the hippocampus matter?, *American Journal of Neuroradiology* 38 (3) (2017) 546–552.
- 495 [22] S. Rathore, M. Habes, M. A. Iftikhar, A. Shacklett, C. Davatzikos, A review on neuroimaging-based classification studies and associated feature extraction methods for alzheimer's disease and its prodromal stages, *NeuroImage* 155 (2017) 530 – 548.
- [23] B. Mwangi, T. S. Tian, J. C. Soares, A review of feature reduction techniques in neuroimaging, *Neuroinformatics* 12 (2) (2014) 229–244.
- 500 [24] E. Westman, A. Simmons, Y. Zhang, J.-S. Muehlboeck, C. Tunnard, Y. Liu, L. Collins, A. Evans, P. Mecocci, B. Vellas, et al., Multivariate analysis of mri data for alzheimer's disease, mild cognitive impairment and healthy controls, *Neuroimage* 54 (2) (2011) 1178–1187.
- 505 [25] C. G. Schwarz, J. L. Gunter, H. J. Wiste, S. A. Przybelski, S. D. Weigand, C. P. Ward, M. L. Senjem, P. Vemuri, M. E. Murray, D. W. Dickson, et al., A large-scale comparison of cortical thickness and volume methods for measuring alzheimer's disease severity, *NeuroImage: Clinical* 11 (2016) 802–812.
- [26] E. E. Bron, M. Smits, W. M. Van Der Flier, H. Vrenken, F. Barkhof, P. Scheltens, J. M. Papma, R. M. Steketee, C. M. Orellana, R. Meijboom,

- 510 et al., Standardized evaluation of algorithms for computer-aided diagnosis
of dementia based on structural mri: the caddementia challenge, *NeuroIm-*
age 111 (2015) 562–579.
- [27] E. Westman, C. Aguilar, J.-S. Muehlboeck, A. Simmons, Regional magnetic
resonance imaging measures for multivariate analysis in alzheimers disease
515 and mild cognitive impairment, *Brain topography* 26 (1) (2013) 9–23.
- [28] R. Cuingnet, E. Gerardin, J. Tessieras, G. Auzias, S. Lehticy, M. O. Habert,
M. Chupin, H. Benali, O. Colliot, Automatic classification of patients with
Alzheimer’s Disease from structural MRI: a comparison of ten methods
using the ADNI database, *Neuroimage* 56 (2) (2011) 766–781.
- 520 [29] A. Lebedev, E. Westman, G. Van Westen, M. Kramberger, A. Lundervold,
D. Aarsland, H. Soininen, I. Kłoszewska, P. Mecocci, M. Tsolaki, et al.,
Random forest ensembles for detection and prediction of alzheimer’s disease
with a good between-cohort robustness, *NeuroImage: Clinical* 6 (2014) 115–
125.
- 525 [30] E. Moradi, A. Pepe, C. Gaser, H. Huttunen, J. Tohka, A. D. N. Initia-
tive, et al., Machine learning framework for early MRI-based Alzheimer’s
conversion prediction in MCI subjects, *Neuroimage* 104 (2015) 398–412.
- [31] J. Tohka, E. Moradi, H. Huttunen, Comparison of feature selection tech-
niques in machine learning for anatomical brain MRI in dementia, *Neu-*
roinformatics (2016) in press.
530
- [32] C. Chu, A.-L. Hsu, K.-H. Chou, P. Bandettini, C. . Lin, A. D. N. Initia-
tive, et al., Does feature selection improve classification accuracy? Impact
of sample size and feature selection on classification using anatomical mag-
netic resonance images, *Neuroimage* 60 (1) (2012) 59–70.
- 535 [33] A. Zandifar, V. Fonov, P. Coupé, J. Pruessner, D. L. Collins, A. D. N.
Initiative, et al., A comparison of accurate automatic hippocampal seg-
mentation methods, *NeuroImage* 155 (2017) 383 – 393.

- [34] J.-L. Chepkoech, K. B. Walhovd, H. Grydeland, A. M. Fjell, Effects of
change in freesurfer version on classification accuracy of patients with
540 alzheimer’s disease and mild cognitive impairment, *Human brain mapping*
37 (5) (2016) 1831–1841.
- [35] O. Voevodskaya, A. Simmons, R. Nordenskjöld, J. Kullberg, H. Ahlström,
L. Lind, L.-O. Wahlund, E.-M. Larsson, E. Westman, A. D. N. Initiative,
et al., The effects of intracranial volume adjustment approaches on multiple
545 regional mri volumes in healthy aging and alzheimer’s disease, *Frontiers in*
aging neuroscience 6.
- [36] Y. Cui, B. Liu, S. Luo, X. Zhen, M. Fan, T. Liu, W. Zhu, M. Park, T. Jiang,
J. S. Jin, et al., Identification of conversion from mild cognitive impairment
to alzheimer’s disease using multivariate predictors, *PloS one* 6 (7) (2011)
550 e21896.
- [37] J. Dukart, M. L. Schroeter, K. Mueller, A. D. N. Initiative, et al., Age
correction in dementia–matching to a healthy brain, *PloS one* 6 (7) (2011)
e22193.
- [38] K. Franke, G. Ziegler, S. Klöppel, C. Gaser, A. D. N. Initiative, et al.,
555 Estimating the age of healthy subjects from t1-weighted MRI scans using
kernel methods: Exploring the influence of various parameters, *Neuroimage*
50 (3) (2010) 883–892.
- [39] J. A. Hanley, B. J. McNeil, The meaning and use of the area under a receiver
operating characteristic (roc) curve., *Radiology* 143 (1) (1982) 29–36.
- 560 [40] C.-C. Chang, C.-J. Lin, LIBSVM: A library for Support Vector Machines,
ACM Trans. Intell. Systems Tech. 2 (2011) 27.
- [41] J. H. Friedman, T. Hastie, R. Tibshirani, Regularization paths for gener-
alized linear models via coordinate descent, *J. Stat. Software* 33 (1) (2010)
1–22.

- 565 [42] T. Hastie, R. Tibshirani, J. Friedman, The elements of statistical learning
- Data Mining, Inference, and Prediction, Second Edition, Springer series
in statistics New York, 2009.
- [43] H. Zou, T. Hastie, Regularization and variable selection via the elastic net,
J. R. Stat. Soc.: Series B 67 (2) (2005) 301–320.
- 570 [44] C. Aguilar, E. Westman, J.-S. Muehlboeck, P. Mecocci, B. Vellas, M. Tso-
laki, I. Kloszewska, H. Soininen, S. Lovestone, C. Spenger, et al., Dif-
ferent multivariate techniques for automated classification of mri data in
alzheimers disease and mild cognitive impairment, Psychiatry Research:
Neuroimaging 212 (2) (2013) 89–98.
- 575 [45] L. Breiman, Random Forests, Machine Learning 45 (1) (2001) 5–32.
- [46] M. N. Wright, A. Ziegler, ranger: A fast implementation of random forests
for high dimensional data in c++ and r, Journal of Statistical Software 77.
- [47] C. Nadeau, Y. Bengio, Inference for the generalization error, in: Advances
in neural information processing systems, 2000, pp. 307–313.
- 580 [48] T. G. Dietterich, Approximate statistical tests for comparing supervised
classification learning algorithms, Neural computation 10 (7) (1998) 1895–
1923.