A Role for Prior Knowledge in Statistical Classification of the Transition from Mild Cognitive Impairment to Alzheimer's Disease

- 6 Initiative¹
- ⁷ ^aDepartment of Statistics, Michigan State University, East Lansing, MI, USA
- ⁸ ^bDepartment of Epidemiology and Biostatistics, College of Human Medicine, Michigan State University, East
- 18 Lansing, MI, USA
- 11
- 12 Handling Associate Editor: Ali Ezzati

Accepted 2 August 2021 Pre-press 25 August 2021

13 Abstract.

- Background: The transition from mild cognitive impairment (MCI) to dementia is of great interest to clinical research on Alzheimer's disease and related dementias. This phenomenon also serves as a valuable data source for quantitative methodological researchers developing new approaches for classification. However, the growth of machine learning (ML)
- approaches for classification may falsely lead many clinical researchers to underestimate the value of logistic regression
- (LR), which often demonstrates classification accuracy equivalent or superior to other ML methods. Further, when faced with many potential features that could be used for classifying the transition, clinical researchers are often unaware of the
- ²⁰ relative value of different approaches for variable selection.
- **Objective:** The present study sought to compare different methods for statistical classification and for automated and theoretically guided feature selection techniques in the context of predicting conversion from MCI to dementia.
- 23 Methods: We used data from the Alzheimer's Disease Neuroimaging Initiative (ADNI) to evaluate different influences of
- automated feature preselection on LR and support vector machine (SVM) classification methods, in classifying conversion
 from MCI to dementia.
- Results: The present findings demonstrate how similar performance can be achieved using user-guided, clinically informed
 pre-selection versus algorithmic feature selection techniques.
- 28 Conclusion: These results show that although SVM and other ML techniques are capable of relatively accurate classification,
- similar or higher accuracy can often be achieved by LR, mitigating SVM's necessity or value for many clinical researchers.
- Keywords: Alzheimer's disease, classification, machine learning, mild cognitive impairment, support vector machine, variable
 selection

*Correspondence to: Andrew R. Bender, Department of Epidemiology and Biostatistics, College of Human Medicine, Michigan State University, East Lansing, MI, USA. E-mail: arbender@msu.edu.

¹Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI)

database (http://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/wpcontent/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf

⁵ Zihuan Liu^a, Tapabrata Maiti^a, Andrew R. Bender^{b,*} and for the Alzheimer's Disease Neuroimaging

32 INTRODUCTION

Alzheimer's disease (AD) is a progressive, age-33 related, neurodegenerative disease and the most 34 common cause of dementia [1–3]. Behaviorally, AD 35 is commonly preceded by mild cognitive impair-36 ment (MCI), a syndrome characterized by declines 37 in memory and other cognitive domains that exceed 38 cognitive decrements associated with normal aging 39 [2, 4]. However, the prodromal symptoms of MCI 40 are not prognostically deterministic: individuals with 41 MCI tend to progress to diagnoses of probable AD 42 at a rate of 8%-15% per year, and many conversions 43 are detectable within 3 years of initial presentation 44 [5–7]. Research efforts to provide new insights into 45 the incidence of MCI-to-AD conversion have focused 46 largely on clinically or biologically relevant features 47 (i.e., neuroimaging markers, clinical exam data, neu-48 ropsychological test scores) and on different methods 49 for statistical classification [8]. 50

For clinical researchers, however, there may be 51 a tendency to conflate more sophisticated, novel 52 analytic approaches and the value of multimodal 53 information from neuroimaging and clinical assess-54 ment. Moreover, whereas statisticians may inherently 55 understand the comparability of different quantita-56 tive approaches, the novelty of both big data and 57 data-driven approaches for studying MCI-to-AD con-58 version may lead clinical researchers to assume that 59 such data-driven methods are inherently superior to 60 more theoretically grounded approaches. Thus, the 61 value of using extant findings and domain exper-62 tise to help guide and constrain the application of 63 newer data-driven approaches capable of capitalizing 64 on emergent big data may be a particularly important 65 consideration for clinical researchers. 66

Statistical classification in clinical research has tra-67 ditionally utilized binary logistic regression (LR). 68 However, key attributes of modern clinical and neu-69 roimaging data, including high dimensionality and 70 the presence of ground truth estimates of pathology 71 and diagnosis provide new opportunities for quantita-72 tive research. This has led to a substantial expansion 73 in the use of data from the Alzheimer's Disease Neu-74 roimaging Initiative (ADNI; http://adni.loni.usc.edu) 75 for quantitative research and methodological devel-76 opment, particularly by researchers utilizing and 77 developing prediction and classification methods in 78 machine learning (ML). Besides LR, support vector 79 machine (SVM) has quickly become the most com-80 mon type of ML classifier for diagnostic prediction 81 and classification with ADNI data. In general, LR 82

works well when the data is linearly separable, and the number of data is greater than the number of features. Moreover, SVM and LR have similar misclassification rates (MCRs) when used to diagnose malignant tumors from imaging data [9, 10].

83

84

85

86

87

88

80

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

Indeed, before the rapid expansion of ML research and applied work over the past decade, many clinical researchers and those outside of engineering and mathematically intensive disciplines had little exposure to classification approaches other than LR. Despite its growing popularity, the relative benefits of SVM or other forms of ML [11, 12] over LR for such classification are not always apparent. Although this may be of little surprise to statisticians and quantitative researchers, such perspectives are often lost on clinical researchers, whose implicit beliefs in the superiority of ML is driven by the volume of publications, rather than through training or empirical demonstration.

Most efforts to develop new classification methods for prediction of MCI-to-AD conversion are well suited to integrate measures from multiple sources such as demographics, clinical rating scores, neuropsychological testing, neuroimaging, genetic markers, etc. However, identifying which combination of features most accurately classifies conversion from MCI to AD is a key challenge for ADNI, and may vary by method. The L_1 norm regularization method (i.e., L_1) is a highly used feature selection technique for LR and SVM. L_1 is popular for addressing circumstances in which the number of features is quite large or even larger than the sample size. Despite some risk of abusing the statistical terminology, the problem is often generically referred to as the "small n, large p" or high dimensional problem. The L_1 technique has dual impacts, namely the algorithm can (i) optimize a higher number of parameters in comparison to sample size, and (ii) reduce the effective number of parameters (i.e., performing variable selection). This powerful technique has been implemented in ADNI data with LR [13]. Furthermore, L_1 and other algorithmic feature selection methods used in ML suffer from one key limitation: they are agnostic to theoretical considerations, and as such, they cannot interpret why selected features are meaningful and important to the model. When sampling from a large pool of features, the algorithmic approaches fail to consider prior knowledge of features and their associations with the relevant systems in variable selection. Therefore, domain expertise and prior knowledge may afford additive or differential value for choosing features and interpreting model

results over algorithmic feature selection methodsalone.

However, most real-world problems occur in the 137 context of additional information about each poten-138 tial feature and its conceptual relationship with the 139 phenomenon being classified. Other than using L_1 140 feature selection, manually trimming the list of 141 potential predictor variables can also protect against 142 over-fitting, and also offers potential insight into 143 why selected features are important to the model. 144 When guided by prior knowledge, user-guided or 145 'manual' feature selection may be a valuable addi-146 tional step to help minimize potentially spurious 147 effects. This perspective is frequently lost on applied 148 researchers, as most commonly used variable selec-149 tion algorithms are context-free-that is, they only 150 look at relationships within the data set, and cannot 151 factor in the wider meanings of variables. Further-152 more, this also means that automated algorithms may 153 identify relationships among a large number of pre-154 dictor variables that are spurious and are unlikely 155 to generalize outside the data set. Although there 156 are a vast number of potential neuroimaging fea-157 tures in ADNI data, the present study focused only 158 on regional brain volumes segmented from struc-159 tural magnetic resonance imaging (MRI) data, the 160 most common neuroimaging datatype for classify-161 ing MCI-to-dementia conversion. In contrast to prior 162 studies that used a limited set of volumetric brain 163 features, the present study utilized data generated 164 by modern multi-atlas segmentation methods and 165 analyses included up to 259 features-anatomically 166 specific gray and white matter volumes. However, the 167 large pool of extant findings from studies evaluat-168 ing regional brain MRI volumetry in prediction and 169 classification of MCI-to-dementia conversion using 170 both limited and expansive feature sets also provides 171 a valuable set of priors for relevant brain regions 172 [14-19]. Thus, applied researchers are often left with 173 the conundrum of more confirmatory approaches that 174 use few regions in classification or more exploratory 175 methods in which prior findings have little value. 176

The present study addressed two questions regard-177 ing commonly used classification approaches for 178 predicting MCI-to-dementia conversion in multi-179 modal data from ADNI. First, we compared 180 performance accuracy of binary LR with SVM in 181 classifying MCI-to-dementia conversion. Second, we 182 asked if applying prior knowledge in feature selection 183 outperforms algorithmic variable selection alone. We 184 hypothesized that 1) LR would perform compara-185 bly to SVM, and 2) user-guided variable selection 186

would outperform algorithmic variable selection alone. This work is intended to demonstrate to clinical researchers the benefit of using ML in an informed fashion, rather than as a 'black box' that obscures clear interpretation. Moreover, we wish to emphasize that this study is not meant to highlight a novel innovation in quantitative methods, but rather to provide an important example to applied researchers regarding the comparable value of ML methods and importance of domain expertise in classification with ADNI data.

MATERIALS AND METHODS

Data used in the preparation of this article were obtained from the ADNI database (http://adni.loni.usc.edu). The ADNI 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 MRI, positron emission tomography (PET), other biological markers, and clinical and neuropsychological assessment can be combined to measure the progression of MCI and early AD. For up-to-date information, see http://www.adni-info.org.

Determination of sensitive and specific markers of preclinical AD and MCI is intended to aid researchers and clinicians to develop new treatments and monitor their effectiveness, as well as reduce the time and cost of clinical trials. Data in the present study came from all sites across the U.S and Canada. All ADNI study participants included in the present analyses were between 55 and 90 years old, spoke English or Spanish as their native language, and had a study partner who provided an independent assessment of functioning.

This study used a subset of the 819 participants from ADNI-1 diagnosed with MCI at baseline and for whom the data from demographic, clinical cognitive assessments, *APOE4* genotyping, and MRI measurements were also available. To evaluate differences in classification performance due to participant inclusion and drop out, we subdivided the sample into two overlapping groups. After applying other criteria for inclusion, Group One included all patients whose follow-up period was at least 36 months (n=265); Group Two consisted of all patients with follow-up assessments at 24 months (n=308). Although the ADNI study protocol includes additional follow-up visits at 6-month intervals, the present study only evaluated baseline data for features (i.e., clinical, 187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

 Table 1

 Sample Sizes by Timing and Diagnosis: Group One and Two

Group	Time	# MCI-S (y=0)	# MCI-C (y=1)	# Total	
One	36 months	101	164	265	
Two	24 months	122	186	308	

Table 1 shows the number of MCI-C, MCI-S, and total subjects in Group One and Two. The number of MCI-C patients is higher than MCI-S patients in both groups.

neuropsychological, brain volumetric) in classifica-236 tion analyses. In addition, identification of stable 237 versus converting clinical outcomes only considered 238 longer-term outcomes based on assessments at 2 and 239 3 years after baseline. The final samples included 240 265 and 308 study participants in Groups One and 241 Two, respectively, who met criteria for inclusion. 242 Both Groups included participants who were stable in 243 their diagnosis (MCI-S) and those who converted to a 244 diagnosis of dementia over the 2 or 3 years (MCI-C). 245 Table 1 shows the participant characteristics. Diag-246 nostic criteria for MCI included a Mini-Mental State 247 Examination (MMSE) score at baseline between 24 248 and 30, a Clinical Dementia Rating (CDR) score of 249 0.5, and a subjective memory complaint, in addition 250 to objective memory loss measured by education-251 adjusted scores on the Logical Memory II subscale 252 of the Wechsler Memory Scale, generally preserved 253 activities of daily living and no dementia. The diag-254 nostic criteria for dementia were an MMSE score 255 between 20 and 26, and a CDR score between 0.5 and 256 1.0. The clinical status of each participant diagnosed 257 with MCI was re-assessed at each follow-up visit and 258 updated to reflect one of several outcomes (e.g., MCI 259 or dementia subtypes). The MCI-C and MCI-S group 260 designations were based on this follow-up clinical 261 diagnosis and marked as either 1 for MCI-C or 0 for 262 MCI-S in classification study. 263

264 Data used in classification

Evaluation of extant reports of common predic-265 tors of conversion from MCI to dementia focused on 266 dimensions of neuropsychological test performance, 267 clinical assessment, genetic data, and regional brain 268 volumes. In the present study, we first divided these 269 variables into two sets of features, with all non-brain 270 volumetric variables in one set and all variables rep-271 resenting regional brain volumes in a second set. In 272 addition, we created a third set of features from the 273 volumetry feature set that only included 26 of the 274 259 brain volumes. Henceforth, we refer to mod-275

els that only include one of these three feature sets as 'single-modality,' whereas models that combine brain and non-brain feature sets are referred to as 'multi-modal.'

Clinical cognitive assessment and genetic data

We considered a total of 19 clinical features as potential predictors of MCI-to-AD progression in our classification analyses. These included the following assessment scores: the MMSE, CDR-Sum of Boxes, Alzheimer's Disease Assessment Scale-cognitive sub-scale (ADAS-cog), Functional Activities Ouestionnaire (FAO) measures of activities of daily living, Trail Making Test-B (TRABSCOR), the immediate and delayed recall components of the Rey Auditory Verbal Learning Test (RAVLT), the Digit-Symbol Coding test (DIGT), and the Digit Symbol Substitution Test from the Preclinical Alzheimer Cognitive Composite (mPACCdigit). We also considered genotype for carriers of the epsilon-4 allele of the apolipoprotein E (APOE) gene [8] as a genetic predictor in this study. Table 2 summarizes all 19 clinical, demographic, and genetic features used in this study. Preliminary comparison of six clinical and genetic predictors by MCI-C and MCI-S subgroups showed five of them (APOE4, ADAS4, CDR, MMSE, and RAVLT.learning) significantly differ between the groups, whereas one (SEX) does not. Figures 1 and 2 illustrate the distribution of these predictors for both groups. Overall, in comparison to MCI-S participants, those in the MCI-C group were more cognitively and functionally impaired at baseline, exhibited greater verbal memory impairments, and included a greater proportion of APOE4 carriers.

MRI data

Structural MRI data were collected according to the ADNI acquisition protocol using T1-weighted scans (GradWarp, B1 Correction, N3, Scaled) [20]. These data included baseline structural MRI scans of 840 ADNI participants, including 230 diagnosed as cognitively normal, 200 with diagnoses of dementia, and 410 diagnosed with MCI. Processing for region-of-interest (ROI)-based volumetric data used in the present study included brain extraction [21] and a multi-atlas, consensus-based label fusion scheme for anatomical parcellation [22] to generate template-based ROIs deformed to individual subject space. MRI scans were automatically segmented into

Characteristics	MCI-S	MCI-C	Test statistic	Р			
Age (y)	74.34 ± 7.78	74.84 ± 6.83	-0.528	$> 0.5^{a}$			
Education (y)	15.57 ± 2.94	15.73 ± 2.91	-0.527	>0.5 ^b			
Sex, % female	33.67%	34.14%	0	1 ^b			
APOE4 carriers %	34.65%	62.19%	17.900	< 0.001 ^a			
CDRSB	1.23 ± 0.61	1.72 ± 0.92	-5.237	< 0.001 ^a			
MMSE score	27.61 ± 1.74	26.82 ± 1.71	3.645	< 0.001 ^a			
ADAS11	8.89 ± 3.79	12.29 ± 4.16	-6.823	< 0.001 ^a			
ADAS13	14.48 ± 5.50	20.01 ± 5.79	-7.795	< 0.001 ^a			
ADASQ4	4.76 ± 2.19	6.77 ± 2.21	-7.339	< 0.001 ^a			
RAVLT.immediate	36.21 ± 10.10	29.10 ± 7.98	6.021	< 0.001 ^a			
RAVLT.learning	4.19 ± 2.74	2.91 ± 2.26	4.231	< 0.001 ^a			
RAVLT.forgetting	4.31 ± 2.59	4.47 ± 2.15	-1.501	< 0.135 ^a			
RAVLT.perc.forgetting	51.55 ± 31.04	72.85 ± 30.45	-5.464	< 0.001 ^a			
LEDLTOTAL	4.96 ± 2.36	3.41 ± 2.66	4.931	< 0.001ª			
DIGTSCOR	40.75 ± 11.09	36.62 ± 10.96	2.882	< 0.005 ^a			
TRABSCOR	109.43 ± 62.94	132.09 ± 71.36	-2.704	< 0.007 ^a			
FAQ	1.50 ± 2.99	4.96 ± 4.79	-7.243	< 0.001 ^a			
mPACCdigit	-5.38 ± 2.96	-8.06 ± 2.96	7.174	< 0.001 ^a			
mPACCtrailsB	-5.47 ± 3.06	-8.22 ± 2.98	7.174	< 0.001 ^a			

Table 2 Clinical Features and Cognitive Assessment Score of Group One

Table only for Group One where has 265 patients and 36 months follow-up time. Values are shown as mean ± standard deviation or percentage. Test statistics and *p*-values for differences between MCI-S and MCI-C are based on (a) *t*-test or (b) chi- square test. MCI-S, non-progressive MCI; MCI-P, progressive MCI; *APOE*, apolipoprotein E; MMSE, Mini-Mental State Examination; RAVLT, The Rey Auditory Verbal Learning Test (immediate: sum of 5 trails; learning: trial 5-trial 1; Forgetting: trial 5-delayed; perc.forgetting: Percent forgetting); DIGT, The Digit- Symbol Coding test; TRAB, Trail Making tests; CDRSB, Clinical Dementia Rating Scaled Response; FAQ, Activities of Daily living Score; ADAS, Alzheimer's Disease Assessment Scale–Cognitive sub-scale; mPACCdigit, the Digit Symbol Substitution Test from the Preclinical Alzheimer Cognitive Composite.



Fig. 1. Comparison of distributions for baseline predictor variables between MCI-S and MCI-C groups. (a) The mean MMSE score in MCI-S is higher than in MCI-C. (b) Mean Learning scores of MCI-C and MCI-S groups are 2.5 and 5.

145 anatomic ROIs spanning the entire brain. An
additional 114 derived ROIs were calculated by combining single ROIs within a tree hierarchy, to obtain
volumetric measurements from larger structures [20].
In total, 259 ROIs were measured and used as potential predictors of MCI-to-dementia progression in this
study.

One of the goals of this study is to investigate if manually selecting predictors improves a model's

331

332

performance. Based on the extant literature [23], we manually selected 26 out of 259 features as theoretically significant predictors of MCI to dementia progression (Table 3) [14–19]. While many brain regions have been reported as showing some relationship to MCI-to-dementia progression, prior reports and reviews clearly implicate hippocampal and entorhinal cortical volumes as markers of such conversion. In addition, we manually selected additional

333

334

335

336

337

338

339

340



Fig. 2. Comparisons between MCI-S and MCI-C groups on baseline predictor variables. The y-axis of panels (a) through (d) represents the number of participants developing AD. Blue and red bars represent non-converters and converters, respectively. Panel (a) shows a greater number of converters than non-converters for both men and women. Panel (b) shows more than half of MCI-C subjects are *APOE4* carriers and approximately 70% MCI-S subjects are non-*APOE4* carriers. Panel (c) shows MCI-S subjects have the relatively lower CDR score and MCI-C subjects have higher CDR score. The number of people in MCI-C group has a downward trend as CDR score increases. Panel (d) shows MCI-C subjects have the relatively higher ADASQ4 score. The average of ASADQ4 score of MCI-S and MCI-C subjects are approximately 5 and 8, respectively.

regions based on their common occurrence across
reports, including cingulate gyrus, precuneus,
amygdala, inferior frontal gyrus, superior parietal
lobule, and lobar white matter volumes.

346 Method and algorithm

In the following section, we utilize binary LR 347 and SVM classification techniques to investigate 348 which approach yields superior discrimination accu-349 racy in the context of ADNI data. Prior comparisons 350 of logistic regression and SVM have reported that 351 SVM requires fewer variables than logistic regres-352 sion to achieve an equivalent level of MCR [10, 353 24]. These also report SVM performs better than 354 LR with microarray expression data [10]. Further-355 more, SVMs have a nice dual form, giving sparse 356 solutions when using the kernel trick. In addition, 357 both methods involve minimizing some cost associ-358 ated with the misclassification based on likelihood 359 ratio for a probabilistic model. Therefore, LR and 360

SVM share common roots in statistical pattern recognition, which we utilize in the comparison of their performance on multi-modal ADNI data.

Logistic regression

LR is the most commonly used machine learning approach for binary classification. In the past decade this has been applied to task of MCI-to-dementia conversion [13, 25, 26]. In the present study, we consider a supervised learning task where we are given M training examples $D = (x_i, y_i), i = 1, ...$ M. Here each $x_i \in \mathbb{R}^N$ is N dimensional feature vectors, and $y_i \in 0, 1$ is a class label. The goal of LR is to model the probability p of a random variable y being 1 or 0 given the experimental data x. The logistic regression model is defined as follows:

$$logit \ p = \log \frac{p}{1-p} \tag{1}$$

Pre-selected MRI features of Group One						
Characteristics	MCI-S	MCI-C	Test statistic	р		
HippoR	3684 ± 438	3366 ± 437	5.735	< 0.001		
HippoL	3414 ± 418	3105 ± 388	5.994	< 0.001		
flWMR	96720 ± 6218	96976 ± 5585	-0.338	0.73		
flWML	93671 ± 5836	94238 ± 5160	-0.802	0.42		
plWMR	47197 ± 3415	47141 ± 3098	0.135	0.89		
plWML	50149 ± 3714	50038 ± 4367	0.242	0.81		
tlWMR	56076 ± 3252	55934 ± 2931	0.359	0.72		
tlWML	55412 ± 3396	55468 ± 3023	-0.136	0.89		
ACgCR	3167 ± 756	3128 ± 641	0.438	0.66		
ACgCL	4104 ± 787	4075 ± 689	0.312	0.76		
EntR	2189 ± 365	1983 ± 373	4.412	< 0.001		
EntL	2050 ± 399	1844 ± 356	4.240	< 0.001		
MCgCR	4176 ± 547	4200 ± 541	-0.341	0.73		
MCgCL	3988 ± 493	4002 ± 559	-0.213	0.83		
MFCR	1581 ± 342	1505 ± 524	1.805	0.07		
MFCL	1566 ± 285	1548 ± 291	0.487	0.62		
OpIFGR	2575 ± 608	2424 ± 546	2.021	0.04		
OpIFGL	2465 ± 550	2361 ± 579	1.466	0.14		
OrIFGR	1252 ± 315	1196 ± 362	1.322	0.18		
OrIFGL	1514 ± 335	1398 ± 356	2.658	< 0.001		
PCgCR	3679 ± 466	3528 ± 415	2.657	< 0.001		
PCgCL	3991 ± 442	3789 ± 424	3.676	< 0.001		
PCuR	10129 ± 1193	9862 ± 1313	1.701	0.09		
PCuL	10005 ± 1263	9759 ± 1299	1.522	0.13		
SPLR	8867 ± 1140	8693 ± 1219	1.180	0.02		
SPLL	8880 ± 1192	8662 ± 1313	1.390	0.17		

Table 3 Pre-selected MRI features of Group One

Values are shown as mean \pm standard deviation or percentage. Test statistics and *p*-values for differences between MCI-C and MCI-S are based on *t*-test. MCI-S, non-progressive MCI; MCI-C, progressive MCI; HippoR, Right Hippocampus; HippoL, Left Hippocampus; flWMR, frontal lobe WM right; flWML, frontal lobe WM left; plWMR, parietal lobe WM right; flWML, temporal lobe WM right; tlWMR, temporal lobe WM right; tlWML, temporal lobe WM left; tlWMR, temporal lobe WM right; tlWML, temporal lobe WM right; tlWML, temporal lobe WM left; tlWMR, temporal lobe WM right; tlWML, temporal lobe WM left; tlWML, temporal lobe WM left; tlWML, temporal lobe WM left; tlWML, temporal lobe WM right; tlWML, temporal lobe WM left; tlWML, temporal lobe with temporal lobe inferior frontal gyrus; OpIFGL, Left OpIFG opercular gyru

Logit, the natural logarithm of the odds, is the key concept that underlies logistic regression. The equation for LR is:

$$\log \frac{P(y_i = 1 | x_i; \beta)}{1 - P(y_i = 1 | x_i; \beta)} = \sum_{j=1}^{N} \beta_j x_{ij}$$
(2)

where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_N)^T$ are the parameters or weights of the logistic regression model, $x_{i,j} = (x_{i1}, \dots, x_{iN})$, $i = 1, \dots, M$. Also, $P(y_i = 1 | x_i; \boldsymbol{\beta})$ is the probability that *ith* MCI patient will develop dementia and $P(y_i = 0 | x_i; \boldsymbol{\beta})$ is the probability that *ith* MCI patient will not develop dementia. Denote $P(y_i = 1 | x_i; \boldsymbol{\beta}) = h(x_i)$,

$$h(x_{i}) = \frac{1}{1 + \exp\left(\sum_{j=1}^{N} -\beta_{j} x_{ij}\right)}$$
(3)

LR is usually trained by minimizing an error function; an appropriate choice of such a function for binary classification problems is the cross-entropy error:

$$e_i(\boldsymbol{\beta}) = -y_i \log(h(x_i)) - (1 - y_i) \log(1 - h(x_i))$$
(4)

The total cost over the data $D = (x^i, y^i), i = 1, ..., M$ is:

$$J(\boldsymbol{\beta}) = -\frac{1}{M} \left[\sum_{i=1}^{M} y_i \log(h(x_i)) - (1 - y_i) \log(1 - h(x_i)) \right]$$
(5)

Consider the problem of finding the maximum likelihood estimate (MLE) of the parameters β for the unregularized logistic regression model. To find the optimized weights β , the total cost needs to be minimized. The optimization function can be written:

$$\boldsymbol{\beta}^{optimal} = min_{\beta} - \frac{1}{M} \left[\sum_{i=1}^{M} y_i log\left(h\left(x_i\right)\right) - (1 - y_i) log\left(1 - h\left(x_i\right)\right) \right]$$
(6)

Solving Equation (6) yields the optimal weights of $\boldsymbol{\beta}$. However, the model-building challenge is to abstract the underlying distribution from the particular instance D of samples because of the relatively small sample size, as compared to the number of features. The problem of replicating the data set instead of identifying the underlying distribution is known as overfitting [27]. To avoid the overfitting problem, it is often necessary to apply a dimension reduction technique. L_1 and L_2 norm are widely used to avoid overfitting, especially when there is a only small number of training examples, or when there is a larger number of features to be learned. L_1 norm or lasso is also often used for feature selection and has been shown to generalize well in the presence of many irrelevant features [28, 29]. L_1 regularization is implemented by adding L_1 norm to the cost function; the cost function and the optimization function were based on the following:

$$J(\boldsymbol{\beta}) = -\frac{1}{M} \left[\sum_{i=1}^{M} y_i log(h(x_i)) - (1-y_i) log(1-h(x_i)) \right] + \lambda |\boldsymbol{\beta}|$$
(7)

and

$$\boldsymbol{\beta}^{optimal} = min_{\beta} \left\{ -\frac{1}{M} \left[\sum_{i=1}^{M} y_i log\left(h\left(x_i\right)\right) - (1 - y_i) log\left(1 - h\left(x_i\right)\right) \right] + \lambda \left|\boldsymbol{\beta}\right| \right\}$$
(8)

365 366 where λ is positive tuning parameter. This Equation (8) is referred to as L_1 regularized logistic regression.

Support vector machine

SVM is another classification and regression method that can handle high dimensional feature vectors. Algorithmically, SVMs build optimal boundaries between data sets by solving a constrained quadratic optimization problem [30–34]. The number of studies applying SVM to evaluate classification of conversion from MCI to dementia has grown over the past decade [1, 2, 5, 8, 23, 35–39].

We briefly review basic support vector machines with linear kernel (SVM-linear) for classification problems: Let $\boldsymbol{\beta}^T h(x) + \boldsymbol{\beta}_0 = 0$ denote an equidistant hyperplane (decision surface) to the closest point of each class on the new space. The goal of SVMs is to find $\boldsymbol{\beta}$ and $\boldsymbol{\beta}_0$ such that $|\boldsymbol{\beta}^T h(x) + \boldsymbol{\beta}_0| = 1$ for all points closer to the hyperplane. In the following classifier construction, one assumes that:

$$\boldsymbol{\beta}^{T} h(x_{i}) + \boldsymbol{\beta}_{0} = \begin{cases} \geq 1 \ if \ y_{i} = 1 \\ \leqslant -1 \ if \ y_{i} = 0 \end{cases}$$
(9)

such that the distance from the closest point of each class to the hyperplane is $1/||\beta||$ and the distance between the two groups is $2/||\beta||$. To maximize the margin, the SVM requires the solution of the following optimization primal problem [40]:

$$\min_{\beta,\beta_0} \sum_{i=1}^{M} \left\{ 1 - y_i \left[\beta_0 + \sum_{j=1}^{N} \beta_j^T h_j \left(x_{ij} \right) \right] \right\}$$
(10)

where h_j is the kernel function which is a linear function for SVM-linear. Specifically, we choose, h_j $(x_j) = x_j$ for *j*-th covariate.

To make the algorithm work for highly correlated features and improve the fitted model's prediction accuracy, we reformulate our optimization by adding L_1 -norm of β , i.e., the *lasso* penalty as follows:

$$\min_{\boldsymbol{\beta},\boldsymbol{\beta}_{0}} \sum_{i=1}^{M} \left\{ 1 - y_{i} \left[\beta_{0} + \sum_{j=1}^{N} \beta_{j}^{T} h_{j} \left(x_{ij} \right) \right] \right\}$$
$$+ \lambda \|\boldsymbol{\beta}\|_{1}$$
(11)

where λ is the tuning parameter that controls the trade-off between loss and penalty. The lasso penalty shrinks the fitted coefficients β towards zero, and hence benefits from the reduction in fitted coefficients' variance.

Experimental design

We built four different classifiers, each designed to classify individual ADNI participants as belonging to

376 377 378

373

374

375

367

385 386

379

380

381

382

383

either the MCI-C group or the MCI-S group: Clas-387 sifier 1 is logistic regression (C-LR); Classifier 2 is 388 logistic regression with L1 norm (C-LR-1); Classifier 389 3 is support vector machine (C-SVM); and Classifier 390 4 is SVM with L_1 norm (C-SVM-1). To test the classi-391 fiers' performance, we constructed five different data 392 sources (Table 4). The first three single-modality data 303 sets included clinical cognitive assessment scores and 394 APOE4 status (CCA), all MRI volumes (ROI-NP), 395 and MRI volumes with preselection (ROI-P), respec-396 tively. Two additional multi-modal data sets were 397 constructed by combining the CCA data separately 398 with ROI-NP and ROI-P data sets (i.e., brain vol-399 umes with and without preselection). Furthermore, it 400 is interesting to note that the number of MCI-S sub-401 jects is 101 (38%) in the Group One and 122 (39%) in 402 Group Two, which makes the data rather imbalanced. 403 Consequently, to precisely report the results obtained 404 from the models, the present study also assessed 405 additional model performance parameters, includ-406 ing AUC score, sensitivity, and specificity (accuracy 407 coefficient is unreliable for imbalanced data). The 408 prediction procedure consisted of three processing 409 stages for Group One (Time = 36 months) and Group 410 Two (Time = 24 months): 1) Split data as training, 411 validation, testing set; 2) Train classifiers using train-412 ing set, tune hyper-parameter using the validation set, 413 and assess classifiers using testing set, then train clas-414 sifiers again using L_1 norm on the same training set; 415 3) Report the testing accuracy, AUC score, sensitivity 416 and specificity of each classifier on single-modality 417 data. Specifically, the first stage used 80% of the sam-418 ple as a training set while the remaining 20% of the 419 data constituted the testing set. In the second stage, 420 the optimal subsets of features of each data source 421 are determined and chosen following application of 422 L_1 norm. We then list the top 10 features of each 423 data set for each of the models. In the last stage, we 424 report AUC score, sensitivity (percent of MCI-C sub-425 jects correctly classified), and specificity (percent of 426 MCI-S subjects correctly classified) as measures of 427 classification accuracy. To protect against over-fitting 428 and to avoid optimistically-biased estimates of model 429 performance, we report 20 measures of predictive 430 performance for each classifier (1-4); for these dif-431 ferent partitions of the data, we report the mean and 432 standard deviation of testing accuracy, AUC score, 433 sensitivity, and specificity (Tables 6 and 7). We also 434 investigate the relationship between the number of 435 features and model performance. Finally, we com-436 pare the performance of LR with SVM based on their 437 ability to handle the problem with a large number 438

Table 4 Modalities

Data sources	# features
Single-modality	
Clinical Cognitive Assessments score and	19
APOE4 data (CCA)	
ROI with no pre-selection data (ROI-NP)	259
ROI with pre-selection data (ROI-P)	26
Multi-modal	
CCA and ROI with no pre-selection data	278
(CCAR-NP)	
CCA and ROI with pre-selection data (CCAR-P)	45

of covariates. Figure 3 illustrates the diagram of the prediction framework.

RESULTS

Cross-validation and choice of λ

We adopted 10-fold cross-validation to tune the hyper-parameters for each model, which included dividing the data into separate sets for training and validation. The ratio of case in training and validation was 8:2. Here, the training set was used to train the model and the validation set was used to select the hyper-parameters. The results of a 10-fold crossvalidation run are summarized with the mean and standard deviation of the model skill scores based on testing data. Cross-validation was also applied to tune the hyper-parameters; λ is used to denote the hyperparameters for both $LR-L_1$ and $SVML_1$. To select the optimized λ , we tried different values of the λ ; results reported here include values of $\lambda = 0.001, 0.01,$ 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, and 0.8 and applied them to the Eq (8) and (11). Next, we selected the λ value based on the best cross-validation score and used the selected λ with Classifiers 2 and 4 to select optimal features. For brevity, the model performance estimates are reported in Tables 6 and 7 for each different modalities, and the top 10 selected features are reported in Table 5. For example, the best λ for ROI-NP- L_1 was 0.01 and the top 3 optimal features selected by LR were left amygdala, right accumbens area, and right middle temporal gyrus. After hyperparameters were selected, we adopted a 10-fold cross-validation again to avoid optimistically-biased estimates of model performance. In each iteration, 212 of the 265 participants are selected by simple random sampling as training cases and the remaining 53 were used as test cases. The approximate 4:1 ratio of training to test cases is, of course, arbitrary.

439 440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473



Fig. 3. Flowchart of the LR and SVM method. A) ROI-P: ROI level data with Pre-selection; B) ROI-NP: ROI level data with No Pre-selection; C) CCAR: Clinical, Cognitive assessments score, *APOE4*, and ROI level data.

475 *Comparison with different modalities*

We compared the performance of each classifier 476 (1-4) on the five different feature sets (Table 4) based 477 on estimates of AUC, sensitivity, and specificity. As 478 shown in Table 6, the results of using LR with L_1 reg-479 ularization (Classifier 2) can achieve the high AUC 480 of 81.2% and sensitivity of 81.4% on single-modality 481 data (CCA), which is considerably better than perfor-482 mance of LR on the other four modalities. Similarly, 483 the best AUC and sensitivity achieved by SVM are 484 81.4% and 81.6% based on the combination of CCA 485 and SVM-L1. Furthermore, we also found the highest 486 accuracy achieved by both classifiers without apply-487

ing regularization is based on the single-modality data (CCA); this indicated both classifiers perform best on single-modality data.

488

489

490

491

Comparison of pre-selection and L₁ norm

We found that using prior knowledge to inform 492 feature selection improves model performance 493 and protects against over-fitting. As shown in 494 Table 6, model performance (i.e., AUC) on ROI-P 495 (64.3%) and CCAR-P (76.3%) outperformed ROI-496 NP (60.6%) and CCAR-NP (60.1%). However, the 497 performance of Classifier 2 on the ROI-NP- L_1 and 498 CCAR-NP-L1 data sets had AUC score of 64.1% 499

Source	LR	-L1 (Classifie	er 2)	SVM-L1 (Classifier 4)			
Data	CCA	ROI-NP	CCAR-NP	CCA	ROI-NP	CCAR-NP	
1	FAQ	AmyL	FAQ	FAQ	AmyL	FAQ	
2	mPACCtrailsB	AccmR	AmyL	Yrs. Educ.	AccmR	AmyL	
3	APOE4	MTGR	ADASQ4	APOE4	AOrGL	AccmR	
4	ADASQ4	HippoL	HippoL	mPACCdigit	PCgGL	AOrGL	
5	Learning	AOrGL	MTGR	ADASQ4	HippoL	PTR	
6	Yrs. Educ.	PrGR	APOE4	Learning	PrGR	AnGR	
7	Forgetting	PCgGL	AOrGL	ADAS11	POrGR	APOE4	
8	mPACCdigit	InfR	Learning	mPACCtrailsB	PTR	PCgGL	
9	ADAS13	POR	mPACCtrailsB	DELTOTAL	LOrGL	Learning	
10	ADAS11	MOGL	mPACCdigit	Forgetting	MOrGL	POrGR	

Table 5Top 10 features of Group One obtained by L_1 regularization

AccmR, Right Accumbens Area; AmyL, Left Amygdala; HippoL, Left Hippocampus; InfR, Right Inf Lat. Vent; AOrGL, Left anterior orbital gyrus; AnGR, Left angular gyrus; LOrGL, Left lateral orbital gyrus; MOGL, Left middle occipital gyrus; MOrGL, Left medial orbital gyrus; MTGR, Right middle temporal gyrus; PCgGL, Left posterior cingulate gyrus; POR, Right parietal operculum; POrGR, Right posterior orbital gyrus; PrGR, Right precentral gyrus; PTR, Right planum temporal.

	Table 6	
LR and SVM performance of Group	p One (Time = 3 years) for models on	single and multi-modal feature sets

Source	LR (Classifier 1 and 2)			SVM (Classifier 3 and 4)				
Modality	Test	Acc% AUC%	Sp%	Sn%	Test	Acc% AUC%	Sp%	Sn% #Features
CCA	$74.3 \pm 6.0,$	$80.8 \pm 7.0,$	62.3 ± 12.1	81.5 ± 6.2	72.4 ± 6.9 ,	80.0 ± 7.3 ,	53.6 ± 13.2	$79.4 \pm 7.7, 19^{(1)}; 19^{(2)}$
ROI-NP	$58.1 \pm 7.0,$	$60.6 \pm 8.1,$	45.5 ± 13.4	65.3 ± 7.9	$59.5 \pm 7.3,$	61.4 ± 8.5 ,	46.5 ± 11.5	$67.3 \pm 8.5, 259^{(1)}; 259^{(2)}$
ROI-P	$64.4 \pm 6.5,$	$64.3 \pm 6.6,$	46.1 ± 10.4	75.0 ± 9.6	$62.1 \pm 5.9,$	64.1 ± 6.2 ,	43.6 ± 9.5	$78.4 \pm 10.4, 26^{(1)}; 26^{(2)}$
CCAR-NP	$57.6 \pm 7.2,$	$60.1 \pm 8.1,$	44.8 ± 12.5	65.1 ± 9.0	$57.8 \pm 6.8,$	$59.1 \pm 7.0,$	45.9 ± 10.4	$65.1 \pm 7.5, 278^{(1)}; 278^{(2)}$
CCAR-P	$72.7 \pm 6.4,$	76.3 ± 6.5 ,	60.5 ± 10.4	80.4 ± 8.2	$66.9 \pm 6.0,$	$69.2 \pm 6.4,$	53.6 ± 13.2	$74.4 \pm 10.5, 45^{(1)}; 45^{(2)}$
$CCA-L_1$	$74.9\pm6.4,$	81.2 ± 6.7 ,	61.3 ± 12.0	83.1 ± 6.6	$72.4 \pm 6.0,$	81.4 ± 6.9 ,	61.6 ± 11.5	$81.6 \pm 5.9, 4^{(1)}; 3^{(2)}$
ROI-NP- L_1	$62.2 \pm 6.6,$	$64.1 \pm 7.9,$	53.1 ± 13.1	68.1 ± 7.2	$62.7 \pm 5.8,$	67.0 ± 6.7 ,	53.7 ± 11.6	$67.7 \pm 7.4, 29^{(1)}; 27^{(2)}$
ROI-P- L_1	$64.4 \pm 6.5,$	$64.3 \pm 6.2,$	46.2 ± 11.0	74.9 ± 9.6	$64.4 \pm 5.7,$	$64.7 \pm 5.8,$	46.7 ± 11.1	$75.4 \pm 8.3, 5^{(1)};17^{(2)}$
CCAR-NP- L_1	$62.6 \pm 7.2,$	$64.0 \pm 8.2,$	51.8 ± 12.7	69.5 ± 7.3	$67.4 \pm 6.4,$	$74.0 \pm 7.4,$	55.7 ± 12.1	$74.1 \pm 7.1, 18^{(1)}; 27^{(2)}$
CCAR-P- L_1	$73.1\pm6.5,$	$77.9\pm5.9,$	61.6 ± 10.5	79.6 ± 7.7	$73.5 \pm 6.2,$	78.5 ± 6.4 ,	61.6 ± 9.3	$80.8 \pm 7.5, 14^{(1)}; 25^{(2)}$

Predictive performance of LR and SVM (mean \pm standard deviation) for all models. Performance estimates include testing accuracy (Test Acc %), area under the cureve (AUC), sensitivity (Sn), and specificity (Sp). The number (#) of features was determined via (1): Classifier 2; (2): Classifier 4.

and 64.0%, while the ROIP- L_1 and CCAR-P- L_1 500 had respective AUC scores of 64.3% and 77.9%; 501 this suggests that user-guided pre-selection signifi-502 cantly improved model performance over L_1 norm. 503 In addition, the SVM (Classifiers 3 & 4) had simi-504 lar and comparable results with LR classifiers. First, 505 as with the LR models, the observed AUC estimates 506 for CCAR-P and ROI-P (69.2% and 64.1%, respec-507 tively), were superior to AUCs from the CCAR-NP 508 (59.1%) and ROI-NP analyses (61.4%). Classifier 4 509 exhibited similar performance on the CCAR-P- L_1 as 510 Classifier 2, with an AUC value of 79.6%—higher 511 than the model for CCAR-NPL₁ (74.0%). Therefore, 512 513 manually selecting features improves model's performance whether L_1 norm is applied, or not. Second, 514 these results show it is necessary and important to 515 use pre-selection because both LR and SVM mod-516 els on CCAR-P- L_1 , with respective AUC estimates 517 of 77.9% and 78.5%, exhibited superior performance 518

over the models without such pre-selection (i.e., LR and SVM on CCAR-NP- L_1 had AUC estimates of 64.0% and 74.0%, respectively).

Comparison of groups one and two

In addition to the results from models of Group One (i.e., MCI-to-AD conversion over 36 months), we also evaluated the performance of Group Two (i.e., MCI-to-AD conversion over 24 months) in an effort to gain further insight regarding possible benefits of shorter or longer assessment periods on classification of the progression of MCI to dementia. Table 7 summarizes the predictive performance of LR and SVM for Group Two. Similarly, we also evaluated classifier performance for single- and multi-modality feature sets. The best result is obtained by using SVM- L_1 model (Classifier 4) on CCAR-P, and its corresponding AUC, Sn and Sp are 76.2%, 60.1%, and

519 520 521

522

523

524

525

526

527

528

529

530

531

532

533

534

		1	1			e			
Source	LR (Classifier 1 and 2)					SVM (Classifier 3 and 4)			
Modality	Test	Acc% AUC%	⊳ Sp%	Sn%	Test 4	Acc% AUC%	Sp%	Sn% #Features	
CCA	$69.9 \pm 5.3,$	76.2 ± 5.5 ,	56.7 ± 9.0	79.3 ± 7.3	$69.4 \pm 5.4,$	75.4 ± 5.5 ,	56.7 ± 8.8	$78.6 \pm 7.1, 19^{(1)}; 19^{(2)}$	
ROI-NP	$58.1 \pm 4.2,$	$58.8 \pm 5.6,$	49.7 ± 7.1	64.4 ± 5.9	$57.8 \pm 5.0,$	$56.6 \pm 6.4,$	50.3 ± 7.1	$62.9 \pm 7.5, 259^{(1)}; 259^{(2)}$	
ROI-P	$63.4 \pm 4.7,$	$65.8 \pm 4.3,$	43.7 ± 10.2	77.8 ± 8.6	$64.5 \pm 4.7,$	$66.2 \pm 5.0,$	44.5 ± 8.5	$79.1 \pm 9.1, 26^{(1)}; 26^{(2)}$	
CCAR-NP	$57.3 \pm 4.0,$	$58.8 \pm 5.4,$	47.5 ± 8.3	64.3 ± 5.8	$56.6 \pm 5.5,$	$56.4 \pm 5.2,$	48.9 ± 7.9	$62.3 \pm 10.4, 278^{(1)}; 278^{(2)}$	
CCAR-P	$70.2 \pm 5.4,$	$74.0 \pm 5.0,$	56.7 ± 9.5	80.6 ± 7.0	$69.5 \pm 4.9,$	72.0 ± 5.3 ,	58.1 ± 8.1	$78.0 \pm 8.2, 45^{(1)}; 45^{(2)}$	
$CCA-L_1$	70.1 ± 4.8 ,	76.3 ± 5.3 ,	56.8 ± 9.9	79.8 ± 7.6	$70.4 \pm 4.9,$	$76.4 \pm 7.7,$	56.8 ± 9.8	$79.4 \pm 7.7, 4^{(1)}; 3^{(2)}$	
ROI-NP- L_1	$62.2 \pm 6.0,$	$64.7 \pm 6.0,$	48.8 ± 9.2	72.0 ± 6.8	$60.8 \pm 4.5,$	$65.9 \pm 6.1,$	53.6 ± 7.5	$64.3 \pm 7.9, 29^{(1)}; 31^{(2)}$	
ROI-P- L_1	$64.1 \pm 4.6,$	$66.8 \pm 3.8,$	42.8 ± 11.3	79.8 ± 8.4	$65.4 \pm 4.0,$	$67.8 \pm 3.9,$	46.3 ± 9.4	$81.8 \pm 7.2, 6^{(1)}; 14^{(2)}$	
CCAR-NP- L_1	$62.6 \pm 6.3,$	$64.8 \pm 6.0,$	49.1 ± 9.1	72.1 ± 6.1	$64.5 \pm 5.1,$	$71.7 \pm 4.8,$	55.4 ± 7.8	$71.4 \pm 8.9, 26^{(1)}; 32^{(2)}$	
CCAR-P- L_1	70.0 ± 5.5 ,	74.3 ± 5.5 ,	57.8 ± 8.0	78.3 ± 8.8	71.3 ± 4.9 ,	76.2 ± 4.7 ,	60.1 ± 7.1	$79.2 \pm 8.5, 14^{(1)}; 27^{(2)}$	

 Table 7

 LR and SVM performance of Group Two (Time = 2 years) for single-data and multi-modal data

For each modality, the predictive performance of LR and SVM are shown (mean \pm standard deviation), including testing accuracy, AUC, sensitivity (Sn), specificity (Sp), # features is the number of features; # features is the number of features; this parameter was determined via (1): Classifier 2; (2): Classifier 4.



Fig. 4. Model performance on ROI feature set by number of features for LR and SVM. Panel (a) shows dramatic growth in AUC with LR as the number of features increases from 1 to 30, and then becoming more static at approximately 74%, i.e., as the number of features increases from 30 to 40, but drops significantly when the number of features reaches to 41. Panel (b) shows the AUC increased dramatically as the number of features grows from 1 to 28, but fluctuated after 29. The optimal number of ROI features for both methods are 29 and 28, and their corresponding optimized AUC were approximately 74.0% and 78.0%.

79.2%, which verifies the assumption that manually 536 selecting techniques improves the model's perfor-537 mance again. However, it warrants mention that all 538 classifiers' performance on the Group One data out-539 performed the same classifiers' performance on the 540 same data sets in Group Two. For example, Classifier 541 2 of Group One on CCA achieved AUC and Sn val-542 ues of 81.2% and 83.1%, which is considerably better 543 than the same classifier of Group Two on CCA (i.e., 544 76.3% and 79.8%). Similarly, Classifier 3 for ROI-545 NP had an AUC of 61.4% for Group One and 56.6% 546 for Group Two. The experimental results indicated 547 superior model performance on data obtained using 548 longer than using shorter follow-up periods. Given 549 the uncertainty in conversion, a longer time window 550

for assessment of cognitive and functional change clearly yields more accurate classification.

Comparison of LR and SVM

In addition to comparing classification between different time windows of assessment, we also compared performance differences between LR and SVM. The results, including models' ability to address the overfitting problem of LR and SVM methods with different modalities are displayed in Tables 6 and 7 and Figs. 4 and 5. First, it is worth noting that both LR and SVM do not work well if no L_1 penalization used, since Classifiers 2 and 4 outperform Classifiers 1 and 3 on the same data set. Second, it

553

554

555

556

557

558

559

560

561

562



Fig. 5. Model performance on CCA feature set by number of features for LR and SVM. Figure (a) shows there is a significant increase in the AUC with LR as the number of features increases from 1 to 5, then there is a slight decrease in the testing accuracy when the number of features is greater than 5. Figure (b) shows the AUC shot up dramatically as the number of features increases from 1 to 4. The optimal number of CCA features obtained by LR and SVM are 5 and 4, and their corresponding optimized AUC are approximately 84.0% and 83.0%.

is worth noting that SVM has a better performance 564 on MRI data when the L1 feature selection method 565 is employed. Third, it was possible to obtain good 566 performance accuracy using LR, which had equiva-567 lent model performance as SVM for "large p" data 568 (ROI-P), as evidenced by respective AUC estimates 569 for Classifiers 1 and 3 of 64.3% and 64.1%. Finally, 570 as shown in Figs. 5 and 4, the SVM method is more 571 stable and robust than LR to the large number of 572 features when n is small. To summarize, the best per-573 formance of Group One was achieved by Classifier 4 574 (SVM with L_1 norm) when using multi-modal, i.e., 575 CCAR-L1, had an AUC of 81.4%. 576

577 DISCUSSION

In this study, we applied two machine learning 578 methods under multiple conditions, to test accuracy 579 in classifying patients with MCI who progress to 580 clinically-defined dementia (MCI-C) from those who 581 remain stable (MCI-S). Using multi-modal data from 582 ADNI, we compared LR and SVM classification 583 accuracy and pre-selection dimensional reduction 584 techniques, i.e., feature selection as informed by prior 585 findings in clinical neuroscience and by L_1 norm. 586 Notably, the present results demonstrate important 587 boundaries for applying feature selection techniques 588 in statistical classification of MCI-to-dementia con-589 version. Specifically, we found that while using L_1 for 590 pre-selection can improve accuracy, it also benefits 591 from a more limited, theoretically based set of feature 592 inputs. In addition, we found that model performance 593 benefited from a longer window of assessment. These 594

results have implications for studies utilizing multimodal data for such classification, including features from clinical neuropsychological assessment, demographic and genetic markers, MRI-based volumetric brain measures, and other modalities.

Comparison of user-defined and L_1 pre-selection for LR and SVM classifiers yielded multiple noteworthy findings, consistent with previously published reports [1, 2, 5, 8, 13, 23, 25, 35]. First, the classification results showed that the model using multi-modal data with cognitive, clinical, and volumetric data (CCAR) achieved better classification accuracy than the methods based on single-modality (CCA, ROI). Moreover, the AUC of CCAR based on LR or SVM was either statistically significantly or at least numerically greater than those based on the single-modality model. Based in AUC, we reported the highest accuracy was observed for CCAR data at 78.5% by L_1 SVM and 77.9% by L1 LR. Second, SVM demonstrated several advantages over LR in discriminating MCI-C from MCI-S (Fig. 4). For one, SVM performance tended to be more stable than LR when the number of features was relatively large. In other words, the model performance of SVM on ROI data remained more stable than LR when using larger numbers of features without user-defined preselection. In particular, SVM performance on ROI data improved as the number of features increased from 20 and 30. In contrast, the AUC values for ROI data sets remained fairly static despite increasing the number of features. However, LR model performance decreased gradually after the number of ROI features reached 40. Third, the classification results clearly

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

demonstrate that manually selecting features on MRI 628 data not only improved the model performance and 629 protected the classifier from overfitting, but also 630 affords easier interpretation of each selected feature's 631 contribution to the model. In addition, we show that 632 pre-selection improves performance: Tables 6 and 7 633 suggest it is the best strategy to obtain the maximum 634 model performance, compared to features selection 635 based on L_1 norm. 636

The present findings can also be interpreted in the 637 context of other reports over the past decade that also 638 investigated the prognostic capacity of brain volume-639 try data to predict the conversion of MCI to dementia, 640 using either SVM or LR, and that also combined 641 volumetry data with other imaging and biomarker 642 modalities such as MRI, functional MRI (fMRI), PET 643 to cerebrospinal fluid (CSF) protein markers [1, 2, 644 5, 8, 13, 23, 25, 35, 41-43]. In addition, one can 645 vary the degrees of non-linearity and flexibility in 646 the model by employing different kernel functions. 647 For example, Young et al (2013) report [8], results 648 from both SVM and Gaussian process (GP) clas-649 sification on MCI progression in ADNI data using 650 MRI, PET, APOE4, and CSF biomarkers. In contrast 651 the present study and with other published work that 652 used MCI-C and MCI-S groups as training and test 653 data sets, they trained a classifier to distinguish cog-654 nitively normal older adults from those diagnosed 655 as probable AD. They reported that the accuracy 656 using GP, an AUC value of 79.5%, was substan-657 tially higher than using any individual modality or 658 using multi-kernel SVM. Other studies of MCI-to-659 dementia classification reporting high accuracy have 660 also implemented other approaches such as multi-661 ple kernel learning (pMKL) classification techniques 662 using clinical, MRI and plasma biomarkers data. One 663 method using this approach to identify the important 664 features first grouped the data set into five differ-665 ent data sources and then applied a filter-wrapper 666 approach of feature selection techniques in combi-667 nation with Joint Mutual Information (JMI) criterion 668 to achieve an AUC of 82% [23]. 669

We also found consistently superior classification 670 performance in patients classified under a longer win-671 dow of assessment. MCI-to-dementia conversion is a 672 process that can take several years to reliably track an 673 individual from onset of amnestic MCI to early-stage 674 dementia [8, 44, 45]. For the modeled features to be 675 of use for classification necessitates well-defined, if 676 not orthogonal classes. However, MCI is not inher-677 ently prodromal to dementia: a large proportion of 678 individuals with MCI never progress, either revert-679

ing to cognitively normal status or remaining rather stable. Furthermore, others may show early evidence of brain atrophy that precedes cognitive impairment by years. In order to account for this variable timing, others have employed methods such as supervised learning using time windows [46]; however, even those methods strongly benefit from longer followup periods. Thus, MCI is an inherently heterogeneous and poorly-defined class, particularly in terms of the relationships between brain characteristics and the likelihood and timing of further cognitive decline.

680

681

682

683

684

685

686

687

688

689

690

The brain volumetric data evaluated in the present 691 study were to limited baseline MRI scans. Alter-692 natively, classifying cognitive decline may benefit 693 from further extending the model to accommo-694 date repeated measurements from longitudinal data. 695 While the inclusion of repeated volumetric data 696 should improve classification accuracy, quantifying 697 the improvements in model performance may also 698 depend on other factors, such as added noise or redun-699 dancy from additional brain parameters, or variability 700 in disease progression. In addition, most recent com-701 putational neuroimaging studies in the past few years 702 have utilized features from multiple neuroimaging 703 modalities [5, 26, 36, 37, 39, 47-50]. For exam-704 ple, when Ding et al. applied SVM with PET and 705 MRI data to classify the transition from MCI to 706 AD, they reported the sensitivity and specificity were 707 66.67% and 64.52% [36]. In addition to PET and 708 structural MRI data, CSF protein markers can be 709 used to predict progression from MCI to AD, in 710 addition to proteomic, demographic, and cognitive 711 data [38, 51, 52]. By applying LR with L_1 norm to 712 CSF markers for classifying individual patients as 713 belonging to either the MCI-C and MCI-S group, 714 one study reported a sensitivity and specificity of 715 80% and 75% [26]. Furthermore, Varatharajah and 716 colleagues (2020) showed SVM-linear outperforms 717 other advanced classification methods, including lin-718 ear classifiers-multiple kernel learning (MKL) with 719 linear kernels, SVM with a linear kernel, and gener-720 alized linear model (GLM), in predicting transition 721 from MCI to AD [42]. In general, LR works well 722 when the data is linearly separable and the number of 723 data is greater than the number of features, whereas 724 SVM with Gaussian Kernel is mostly used when 725 the data is not linearly separable. In addition to LR 726 and SVM, deep neural network approaches also offer 727 benefits [41, 53], but have not had the extent of appli-728 cation in ADNI data as SVM and LR. Using a novel 729 LR, artificial neural network (ANN) model and deci-730 sion tree (DT) model for classifying the progression 731 of MCI to AD, Kuang (2021) reported that the ANN
exhibited the highest sensitivity at 82.1% [43].

In conclusion, models applying prior knowledge 734 for classification and prediction of MCI-to-dementia 735 conversion outperform those without pre-selection. 736 This theoretically guided pre-selection of features 737 from MRI-based regional brain volumes appears to 738 protect the model against over-fitting. In addition, the 739 present findings demonstrate that SVM classifier per-740 formance is more stable than LR for dealing with the 741 "large p" problem. Clinical researchers should both 742 note the value of evaluating different classification 743 and pre-selection approaches in application to clini-744 cal or research questions and be mindful that not all 745 machine learning techniques are equally beneficial 746 for modeling specific clinical outcomes. 747

748 ACKNOWLEDGMENTS

751

752

The research is partially supported by NSF-DMS1945824 and 1924724.

We are grateful to the patients and their families who participated in the ADNI.

Data collection and sharing for this project was 753 funded by the Alzheimer's Disease Neuroimag-754 ing Initiative (ADNI) (National Institutes of Health 755 Grant U01 AG024904) and DOD ADNI (Department 756 of Defense award number W81XWH-12-2-0012). 757 ADNI is funded by the National Institute on Aging, 758 the National Institute of Biomedical Imaging and 759 Bioengineering, and through generous contributions 760 from the following: AbbVie, Alzheimer's Asso-761 ciation; Alzheimer's Drug Discovery Foundation; 762 Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-763 Myers Squibb Company; CereSpir, Inc.; Cogstate; 764 Eisai Inc.; Elan Pharmaceuticals, Inc.; Eli Lilly and 765 Company; EuroImmun; F. Hoffmann-La Roche Ltd 766 and its affiliated company Genentech, Inc.; Fujire-767 bio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer 768 Immunotherapy Research & Development, LLC.; 769 Johnson & Johnson Pharmaceutical Research & 770 Development LLC.; Lumosity; Lundbeck; Merck 771 & Co., Inc.; Meso Scale Diagnostics, LLC.; Neu-772 roRx Research; Neurotrack Technologies; Novartis 773 Pharmaceuticals Corporation; Pfizer Inc.; Piramal 774 Imaging; Servier; Takeda Pharmaceutical Company; 775 and Transition Therapeutics. The Canadian Institutes 776 of Health Research is providing funds to support 777 ADNI clinical sites in Canada. Private sector con-778 tributions are facilitated by the Foundation for the 779 National Institutes of Health (http://www.fnih.org). 780

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

Authors' disclosures available online (https:// www.j-alz.com/manuscript-disclosures/20-1398r3).

REFERENCES

- [1] Zhang D, Shen D (2012) Multi-modal multi-task learning for joint prediction of multiple regression and classification variables in Alzheimer's disease. *Neuroimage* **59**, 895-907.
- [2] Zhang D, Shen D (2012) Predicting future clinical changes of MCI patients using longitudinal and multimodal biomarkers. *PLoS One* 7, 1-15.
- [3] Korolev IO (2014) Alzheimer's disease: A clinical and basic science review. *Med Stud Res J* **4**, 24-33.
- [4] Petersen RC, Roberts RO, Knopman DS, Boeve BF, Geda YE, Ivnik RJ (2009) Mild cognitive impairment: Ten years later. Arch Neurol 66, 1447-1455.
- [5] Cui Y, Liu B, Luo S, Zhen X, Fan M, Liu T (2011) Identification of conversion from mild cognitive impairment to Alzheimer's disease using multivariate predictors. *PLoS One* 6, e21896.
- [6] Farlow MR (2009) Treatment of mild cognitive impairment (MCI). Curr Alzheimer Res 6, 262-267.
- [7] Salazar D, V´elez J, Salazar J (2012) A relationship between the transient structure in the monomeric state and the aggregation propensities of alpha-synuclein and beta-synuclein. *Biochemistry* 53, 7170-7183.
- [8] Young J, Modat M, Cardoso MJ, Mendelson A, Cash D, Ourselin S (2013) Accurate multimodal probabilistic prediction of conversion to Alzheimer's disease in patients with mild cognitive impairment. *Neuroimage Clin* 2, 735-745.
- [9] Chen S, Hsiao Y, Huang Y (2009) Comparative analysis of logistic regression, support vector machine and artificial neural network for the differential diagnosis of benign and malignant solid breast tumors by the use of threedimensional power Doppler imaging. *Korean J Radiol* 10, 464-471.
- [10] Salazar D, Velez J, Salazar J (2012) Comparison between SVM and logistic regression: Which one is better to discriminate. *Rev Colomb Estad* 35, 223-237.
- [11] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E (2011) Scikit-learn: Machine learning in Python. J Mach Learn Res 12, 2825-2830.
- [12] McKinney W (2010) Data structures for statistical computing in Python. Proceedings of the 9th Python in Science Conference, pp. 51-56.
- [13] Ye J, Farnum M, Yang E, Verbeeck R, Lobanov V, Raghavan N (2012) Sparse learning and stability selection for predicting MCI to AD conversion using baseline ADNI data. *BMC Neurol* 12, 46.
- [14] Chapman RM, Mapstone M, McCrary JW, Gardner MN, Porsteinsson A, Sandoval TC, Guillily MD, Degrush ED, Reilly LA (2011) Predicting conversion from mild cognitive impairment to Alzheimer's disease using

781

782

783

784

785

786

787

788

789

790

701

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

neuropsychological tests and multivariate methods. J Clin Exp Neuropsychol 32, 187-189.

- Ewers M, Walsh C, Trojanowski JQ, Shaw LM, Petersen 843 [15] RC, Jack CR, Feldman HH, Bokde AL, Alexander GE, 844 Scheltens P, Vellas B, Dubois B, Weiner M, Hampel 845 H (2012) Prediction of conversion from mild cognitive 846 impairment to Alzheimer's disease dementia based upon 847 biomarkers and neuropsychological test performance. Neu-848 849 robiol Aging 33, 1203-1214.
 - [16] Tabatabaei-Jafari H, Shaw ME, Cherbuin N (2015) Cerebral atrophy in mild cognitive impairment: A systematic review with meta-analysis. Alzheimers Dement (Amst) 1, 487-504.
 - [17] Misra C, Fan Y, Davatzikos C (2009) Baseline and longitudinal patterns of brain atrophy in MCI patients, and their use in prediction of short-term conversion to AD: Results from ADNI. Neuroimage 44, 1415-1422.
 - [18] Eckerstr"om C, Olsson E, Borga M, Ekholm S, Ribbelin S, Rolstad S, Starck G, Edman A, Wallin A, Malmgren H (2008) Small baseline volume of left hippocampus is associated with subsequent conversion of MCI into dementia: The Goteborg MCI study. J Neurol Sci 271, 48-59.
 - [19] Risacher SL, Saykin AJ, West JD, Shen L, Firpi HA, McDonald BC (2009) Baseline MRI predictors of conversion from MCI to probable AD in the ADNI cohort. Curr Alzheimer Res 6, 347-361.
 - [20] DoshivJ, Erus G, Rozycki M, Davatzikos C (2016) Hierarchical parcellation of MRI using multi-atlas labeling methods. Alzheimer's Disease Neuroimaging Initiative, Philadelphia, PA
- [21] Doshi J, Erus G, Ou Y, Gaonkar B, Davatzikos C (2013) 870 Multi-atlas skull-stripping. Acad Radiol 20, 1566-1576. 871
- 872 [22] Doshi J, Erus G, Ou Y, Resnick S, Gur R, Satterthwaite T, Davatzikos C (2015) MUSE: Multiatlas region Segmentation utilizing Ensembles of registration algorithms and parameters, and locally optimal atlas selection. Neuroimage 127, 186-195.
 - Korolev I, Symonds L, Bozoki A (2016) Predicting pro-[23] gression from mild cognitive impairment to Alzheimer's dementia using clinical, MRI, and plasma biomarkers via probabilistic pattern classification. PLoS One 11, 2.
 - Verplancke T, Van Looy S, Benoit D, Vansteelandt S, [24] Depuydt P, De Turck F, Decruyenaere J (2009) Support vector machine versus logistic regression modeling for prediction of hospital mortality in critically ill patients with haematological malignancies. BMC Med Inform Decis Mak 8, 56.
 - Devanand DP, Liu X, Tabert MH, Pradhaban G, Cuasay K, [25] Bell K (2008) Combining early markers strongly predicts conversion from mild cognitive impairment to Alzheimer's disease. Biol Psychiatry 64, 871-879.
 - Llano DA, Bundela S, Mudar RA, Devanarayan V (2017) A [26] multivariate predictive modeling approach reveals a novel CSF peptide signature for both Alzheimer's Disease state classification and for predicting future disease progression. PLoS One 12, 1-18.
- [27] Stephan D, Lucila O (2008) Logistic regression and artifi-897 cial neural network classification models: A methodology review. BMC Med Inform Decis Mak 8, 56-64.
- Lee S, Lee H, Abbeel P, Ng Andrew (2006) Efficient L1 regu-[28] 899 larized logistic regression. Association for the Advancement 900 of Artificial Intelligence, pp. 401-408. 901
 - [29] Tibshirani R (1996) Regression shrinkage and selection via the lasso. J R Stat Soc Series B 58, 267-288.
- [30] Cristianini N, Shawe-Taylor J (2000) An introduction to 904 support vector machines and other kernelbased learning 905

methodsr. Cambridge University Press, Cambridge, United Kingdom.

- [31] Scholkopf B, Smola A (2002) Learning with kernels: Support vector machines, regularization, optimization, and beyond. MIT Press, Boston.
- [32] Vapnik V (1998) Statistical learning theory. John Wiley, New York.
- Vapnik V (1995) The nature of statistical learning theory. [33] Springer-Verlag, New York.
- [34] Vapnik V (1998) The support vector method of function estimation. Kluwer Academic Publisher, Boston, pp. 267-288
- [35] Hinrichs C, Singh V, Xu G, Johnson SC (2011) Predictive markers for AD in a multi-modal framework: An analysis of MCI progression in the ADNI population. Neuroimage 55, 574-589.
- Ding J, Huang Q (2017) Prediction of MCI to AD con-[36] version using Laplace Eigenmaps learned from FDG and MRI images of AD patients and healthy controls. 2nd International Conference on Image, Vision and Computing (ICIVC), pp. 660-664,
- [37] Hojjati S, Ebrahimzadeh A, Khazaee A, Babajani-Feremi A (2017) Predicting conversion from MCI to AD using restingstate fMRI, graph theoretical approach and SVM. JNeurosci Methods 282, 69-80.
- [38] Davatzikos C, Bhatt P, Shaw L, Batmanghelich K, Trojanowski J (2011) Prediction of MCI to AD conversion, via MRI, CSF biomarkers, and pattern classification. Neurobiol Aging 32, 19-27.
- [39] Wei R, Li C, Fogelson N, Li L (2016) Prediction of conversion from mild cognitive impairment to Alzheimer's disease using MRI and structural network features. Front Aging Neurosci 8, 76.
- [40] Zhu J, Rosset S, Hastie T, Tibshirani R (2003) 1-Norm Support Vector Machines. MIT Press, pp. 49-56.
- [41] Li F, Tran L, Thung K, Ji S, Shen D, Li J (2015) A robust deep model for improved classification of AD/MCI patients. IEEE J Biomed Health Inform 19, 1610-1616.
- [42] Varatharajah Y, Ramanan VK, Iyer R, Vemuri P (2019) Predicting short-term MCI-to-AD progression using imaging, CSF, genetic factors, cognitive resilience, and demographics. Sci Rep 2235, 9.
- [43] Kuang J, Zhang P, Cai T, Zou Z, Li L, Wang N, Wu L (2021) Prediction of transition from mild cognitive impairment to Alzheimer's disease based on a logistic regression-artificial neural network decision tree model. Geriatr Gerontol Int 21, 43-47.
- Mitchell A, Shiri-Feshki M (2008) Temporal trends in the [44] long term risk of progression of mild cognitive impairment: A pooled analysis. J Neurol Neurosurg Psychiatry 79, 1386-91.
- [45] Lee S, Bachman A, Yu D, Lim J, Ardekani B (2016) Predicting progression from mild cognitive impairment to Alzheimer's disease using longitudinal callosal atrophy. Alzheimers Dement (Amst) 2, 68-74.
- [46] Pereira T, Lemos L, Cardoso S, Silva D, Rodrigues A, Santana I, DeMendon, ca A, Guerreiro M, Madeira SC (2017) Predicting progression of mild cognitive impairment to dementia using neuropsychological data: A supervised learning approach using time windows. BMC Med Inform Decis Mak 17, 110.
- Shen T, Jiang J, Li Y, Wu P, Zuo C, Yan Z (2018) A mul-[47] tivariate predictive modeling approach reveals a novel CSF peptide signature for both Alzheimer's Disease state classification and for predicting future disease progression. 40th

906

841

842

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

802

893

894

895

896

898

902

17

Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), pp. 738-741.

[48] Menikdiwela M, Nguyen C, Shaw M (2018) Deep learning
 on brain cortical thickness data for disease classification.
 Digital Image Computing: Techniques and Applications (*DICTA*), pp. 1-15.

971

972

- [49] Minhas S, Khanum A, Riaz F, Alvi A, Khan SA (2017)
 A nonparametric approach for mild cognitive impairment to AD conversion prediction: Results on longitudinal data.
 IEEE J Biomed Health Inform 21, 1403-1410.
- [50] Wang B, Hong R, Xu Y, Zhou F, Wang P (2016) Identifying
 mild cognitive impairment conversion to Alzheimer's disease from medical image information. *IEEE International Conference on Consumer Electronics-Taiwan (ICCE-TW)*,
 pp. 1-2.
- [51] Shaffer J, Petrella J, Sheldon F, Choudhury K, Calhoun V and Coleman R, Doraiswamy P (2012) Predicting cognitive decline in subjects at risk for Alzheimer disease by using combined cerebrospinal fluid, MR imaging, and PET biomarkers. *Radiology* 266, 583-591.
- [52] Cheng B, Zhang D, Shen D (2012) Domain transfer learning for MCI conversion prediction. *Med Image Comput Comput Assist Interv* 15, 82-90.
- [53] Suk H. I, Shen D (2013) Deep learning-based feature representation for AD/MCI classification. *Med Image Comput* 6, 583-590.