

1 An Efficient Radius-incorporated MKL Algorithm for  
2 Alzheimer's Disease Prediction

3 Xinwang Liu<sup>a</sup>, Luping Zhou<sup>b</sup>, Lei Wang<sup>b</sup>, Jian Zhang<sup>c</sup>, Jianping Yin<sup>a</sup>,  
4 Dinggang Shen<sup>d</sup>

5 <sup>a</sup>*School of Computer, National University of Defense Technology, Changsha, China,*  
6 *410073*

7 <sup>b</sup>*School of Computer Science and Software Engineering, University of Wollongong,*  
8 *NSW, Australia, 2522.*

9 <sup>c</sup>*Faculty of Engineering and Information Technology, University of Technology, Sydney,*  
10 *NSW, Australia, 2007.*

11 <sup>d</sup>*The Department of Radiology and Biomedical Research Imaging Center (BRIC),*  
12 *University of North Carolina, Chapel Hill, NC 27599 USA.*

---

13 **Abstract**

Integrating multi-source information has recently shown promising performance in predicting Alzheimer's disease (AD). Multiple kernel learning (MKL) plays an important role in this regard by learning the combination weights of a set of base kernels via the principle of margin maximisation. The latest research on MKL further incorporates the radius of minimum enclosing ball (MEB) of training data to improve the kernel learning performance. However, we observe that directly applying these radius-incorporated MKL algorithms to AD prediction tasks does not necessarily improve, and sometimes even deteriorate, the prediction accuracy. In this paper, we propose an improved radius-incorporated MKL algorithm for AD prediction. First, we redesign the objective function by approximating the radius of MEB with its upper bound, a linear function of the kernel weights. This approximation makes the resulting optimisation problem convex and globally solvable.

Second, instead of using cross-validation, we model the regularisation parameter  $C$  of the SVM classifier as an extra kernel weight and automatically tune it in MKL. Third, we theoretically show that our algorithm can be reformulated into a similar form of the SimpleMKL algorithm and conveniently solved by the off-the-shelf packages. We discuss the factors that contribute to the improved performance and apply our algorithm to discriminate different clinic groups from the benchmark ADNI data set. As experimentally demonstrated, our algorithm can better utilise the radius information and achieve higher prediction accuracy than the comparable MKL methods in the literature. In addition, our algorithm demonstrates the highest computational efficiency among all the comparable methods.

14 *Keywords:* Multiple Kernel Learning, Radius-margin Bound, Support  
15 Vector Machines, Alzheimer’s disease, Neuroimaging

---

## 16 **1. Introduction**

17 Pattern recognition techniques have been extensively applied to the anal-  
18 ysis and diagnosis of medical diseases, and their effectiveness and significance  
19 have been well demonstrated in the literature [1, 2, 3]. In particular, accu-  
20 rate classification of people contracting a disease and the healthy population  
21 helps not only treatment but also early prevention. Therefore, developing  
22 better classification methods in this regard is highly desired. In this paper,  
23 we aim to develop a new pattern classification algorithm that can achieve  
24 improved classification performance when applied to Alzheimer’s disease.

25 Alzheimer’s disease (AD in short) is the most common neurodegenerative  
26 disease, covering 60% ~ 70% age-related dementia [4]. It is a fatal disease

27 that worsens as it progresses. Mild cognitive impairment (MCI) is a precursor  
28 of AD. It is heterogeneous, with a conversion rate of 15% per year to  
29 AD [5]. Considering the immense cost on looking after AD patients, early  
30 identification of MCI and AD patients is of great significance. As a result,  
31 the following two classification tasks become important: i) discriminating  
32 MCI patients from the healthy population; and ii) discriminating the MCI  
33 patients who will convert to AD from those who will not. Since the two tasks  
34 can generally be viewed as predicting whether a person will develop towards  
35 or into AD, we call them collectively “AD prediction” for short in this paper.

36 Recent studies have demonstrated neuroimaging techniques as an important  
37 meanings for AD analysis [6, 7]. For example, magnetic resonance  
38 imaging (MRI) shows grey matter morphometry, and Fluorodeoxyglucose  
39 (FDG) positron emission tomography (PET) shows metabolic activity. In  
40 this case, more effective AD prediction methods have been developed by  
41 combining the complementary information carried by these imaging modalities  
42 [8, 9]. As seen in the recent literature, the combination methods can be  
43 performed at feature level [10, 11, 8, 12, 13, 14] or classifier level [15]. A common  
44 practice of feature-level combination is to concatenate the features from  
45 different modalities into a long feature vector [10, 11] and use it for classification.  
46 However, such concatenation usually requires proper normalisation  
47 of the features from different modalities. Otherwise, classification could be  
48 dominated by the features that have large variation but are not necessarily  
49 discriminative, leading to less satisfying classification performance.

50 In the past several years, multiple kernel learning (MKL) has shown superior  
51 performance to the methods using feature-level combination on AD

52 prediction [15]. MKL is an important extension of support vector machines  
53 (SVM) [16] for handling multiple information sources. By predefining one  
54 (or multiple in general) “base” kernel function for each source, MKL aims to  
55 find the optimal linear combination weights of these kernels by maximising  
56 classification-performance-related criteria such as the margin of two classes.  
57 One of the representative algorithms is SimpleMKL [17]. It has been used  
58 for AD prediction by combining multiple modalities such as MRI, PET, and  
59 cerebrospinal fluid (CSF) parameters [12, 13, 14]. Due to its promising clas-  
60 sification performance and solid theoretical foundation, SimpleMKL [17] is  
61 regarded as the state-of-the-art for AD prediction with multiple modalities.

62 Recent research [18, 19, 20] proposes to use more sophisticated criteria to  
63 optimise the kernel weights. In addition to the margin of two classes, these  
64 criteria consider the radius of minimum enclosing ball (MEB) of training  
65 data. The logic lies at that the radius affects the generalisation performance  
66 of SVM and it varies with the kernel weights. Hence, this radius shall be  
67 considered when seeking the optimal weight values. In the following, we call  
68 the MKL algorithms in [18, 19, 20] “radius-incorporated MKL” for short.

69 Our study observes that when applied to AD prediction, these radius-  
70 incorporated MKL algorithms do not necessarily improve, sometimes even  
71 deteriorate, the classification performance. By looking into this, we find  
72 that the tasks of AD prediction often have a small number of training sam-  
73 ples and the involved classes are usually difficult to differentiate. Based on  
74 this observation, we hypothesise that the following two issues lead to the  
75 unsatisfying classification performance of these radius-incorporated MKL al-  
76 gorithms: i) their objective functions are not convex. This usually leads to

77 a locally (rather than globally) optimal solution. Unless the locally optimal  
78 solution is close enough to the (unknown in practice) globally optimal solu-  
79 tion, the kernel weights will not be properly optimised; ii) Essentially as an  
80 SVM classifier, MKL also needs to tune the regularisation parameter  $C$  to  
81 attain good classification. The above radius-incorporated MKL algorithms  
82 employ multi-fold cross-validation technique to tune  $C$ . Nevertheless, when  
83 the number of training samples is small, this technique will become less reli-  
84 able and may select an inappropriate value for  $C$ . Such a selection could lead  
85 to poor classification performance, especially when the classes are difficult to  
86 separate, as in the tasks of AD prediction.

87 To address the above two issues, we propose an improved radius-incorporated  
88 MKL algorithm. Firstly, to address the issue of non-convexity, we employ an  
89 approximation to the radius of MEB in our objective function, rather than  
90 directly using the radius as existing algorithms. This approximation can be  
91 shown as a linear function of the kernel weights. This makes our objective  
92 function convex and a globally optimal solution is therefore guaranteed, as  
93 proved in this paper. Also, we discuss the relationship between this approx-  
94 imation and the original radius to give a theoretical support for using this  
95 approximation. Secondly, to address the issue of tuning the parameter  $C$ , we  
96 do not use cross-validation. Instead, we define an extra dummy base kernel  
97 and relate  $C$  to the weight of this kernel. In doing so,  $C$  can be tuned with  
98 the other weights in MKL, and this mitigates the reliability issue of cross-  
99 validation in the case of small sample. This trick of tuning  $C$  has been used  
100 for model selection of SVM [21, 22] and kernel learning [23]. However, it has  
101 not been integrated into the radius-incorporated MKL algorithms, and we

102 find that tuning  $C$  in this way can effectively help improving the classification  
103 performance of MKL on the tasks of AD prediction.

104 In addition, the radius-incorporated MKL algorithm proposed in this pa-  
105 per brings computational advantage. Our objective function can be trans-  
106 formed into a form similar to that in the SimpleMKL. This allows our  
107 algorithm to be readily implemented by existing software packages. This  
108 merit does not apply to existing radius-incorporated MKL algorithms, which  
109 need more sophisticated optimisation algorithms. Also, as mentioned above,  
110 our algorithm tunes the parameter  $C$  via optimisation instead of timing-  
111 consuming cross-validation. These factors contribute to the higher compu-  
112 tational efficiency, which will be experimentally demonstrated.

113 Experimental studies are conducted on 11 UCI machine learning bench-  
114 mark data sets and three AD prediction tasks. We compare our algorithm  
115 with a set of state-of-the-art MKL algorithms, including unweighted average  
116 MKL, SimpleMKL [17], radius-incorporated algorithms in [18, 19, 20], and  
117 non-sparse MKL algorithms [24]. As will be demonstrated, our algorithm can  
118 achieve better classification performance on AD prediction tasks and higher  
119 computational efficiency than existing algorithms in comparison.

120 The rest of this paper is organised as follows. The background on MKL  
121 is reviewed in Section 2. In Section 3, we develop our algorithm and analyse  
122 its properties. After that, two factors contributing to the improvement of  
123 our algorithm are discussed, and two additional algorithms are designed to  
124 experimentally verify this discussion. Section 5 reports our experimental  
125 study and the conclusion is drawn in Section 6.

126 **2. Background**

127 Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a set of  $n$  training samples.  $\mathbf{x}_i$  ( $i = 1, \dots, n$ ) is a  $d$ -  
 128 dimensional column vector in a Euclidean space  $\mathbb{R}^d$ . Let  $y_i$  be the class label  
 129 of  $\mathbf{x}_i$ , and its value is  $+1$  or  $-1$ . Let  $\phi(\cdot) : \mathbb{R}^d \rightarrow \mathcal{H}$  be a probably nonlinear  
 130 mapping from  $\mathbb{R}^d$  to a higher-dimensional feature space  $\mathcal{H}$ . A kernel function  
 131 of  $\mathbf{x}_i$  and  $\mathbf{x}_j$  is defined as the inner product between  $\phi(\mathbf{x}_i)$  and  $\phi(\mathbf{x}_j)$  [25]. It  
 132 is expressed as  $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ , where  $\top$  denotes the transpose of a  
 133 vector. The mapping  $\phi(\cdot)$  usually cannot be explicitly computed.

134 *2.1. Traditional MKL*

135 As previously mentioned, MKL employs a set of base kernels. Let  $m$  be  
 136 the number of base kernels, and the  $p$ -th kernel is denoted by  $k_p(\cdot, \cdot)$ <sup>1</sup>, where  
 137  $p = 1, \dots, m$ . Accordingly, let  $\phi_p(\cdot) : \mathbb{R}^d \rightarrow \mathcal{H}_p$  be a probably nonlinear  
 138 mapping associated with the  $p$ -th base kernel, where  $\mathcal{H}_p$  is the  $p$ -th feature  
 139 space. It is known that  $k_p(\cdot, \cdot) = \phi_p(\cdot)^\top \phi_p(\cdot)$  by definition. Let  $\gamma_p$  be the  
 140 weight of  $k_p(\cdot, \cdot)$  and let  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_m)^\top$  be an  $m$ -dimensional column  
 141 vector. Let  $k(\cdot, \cdot; \boldsymbol{\gamma})$  denote a linear combination of these base kernels and  
 142 it is expressed as  $k(\cdot, \cdot; \boldsymbol{\gamma}) = \sum_{p=1}^m \gamma_p k_p(\cdot, \cdot)$ . Note that we explicitly show  
 143  $\boldsymbol{\gamma}$  as a parameter of this kernel to emphasise its dependence on these kernel  
 144 weights. In this paper we call  $k(\cdot, \cdot; \boldsymbol{\gamma})$  the “final” kernel for short. Defining  
 145  $\phi(\cdot; \boldsymbol{\gamma}) = [\sqrt{\gamma_1} \phi_1(\cdot)^\top, \dots, \sqrt{\gamma_m} \phi_m(\cdot)^\top]^\top$ , it is not difficult to see that

146 
$$k(\cdot, \cdot; \boldsymbol{\gamma}) = \sum_{p=1}^m \gamma_p k_p(\cdot, \cdot) = \sum_{p=1}^m (\sqrt{\gamma_p} \phi_p(\cdot))^\top (\sqrt{\gamma_p} \phi_p(\cdot)) = \phi(\cdot; \boldsymbol{\gamma})^\top \phi(\cdot; \boldsymbol{\gamma}). \quad (1)$$

---

<sup>1</sup>Note that the symbol  $k_p(\cdot, \cdot)$  is used to emphasise the kernel “function”, while  $k_p(\mathbf{x}_i, \mathbf{x}_j)$  is used to emphasise the kernel function “value”.

147 This result shows that *conceptually, MKL can be viewed as mapping a sample*  
 148  *$\mathbf{x}$  onto a feature space  $\mathcal{H}(\boldsymbol{\gamma})$  via  $\phi(\cdot; \boldsymbol{\gamma})$  and using a single kernel  $k(\cdot, \cdot; \boldsymbol{\gamma})$ .*  
 149 This concept will be frequently used to derive our algorithm.

150 As mentioned in Section 1, MKL aims to seek the optimal kernel weights  
 151  $\boldsymbol{\gamma}$ . Most of existing MKL algorithms [17, 24, 26] find the optimal  $\boldsymbol{\gamma}$  by  
 152 maximising the margin of two classes as

$$\begin{aligned}
 & \min_{\boldsymbol{\gamma}, \boldsymbol{\omega}, b, \boldsymbol{\xi}} \frac{1}{2} \|\boldsymbol{\omega}\|^2 + C \sum_{i=1}^n \xi_i \\
 & s.t. \ y_i(\boldsymbol{\omega}^\top \phi(\mathbf{x}_i; \boldsymbol{\gamma}) + b) \geq 1 - \xi_i, \ \xi_i \geq 0, \forall i; \ \sum_{p=1}^m \gamma_p = 1, \ \gamma_p \geq 0, \ \forall p.
 \end{aligned} \tag{2}$$

154 where  $\boldsymbol{\omega}$  is the normal of the SVM separating hyperplane,  $b$  the bias of  
 155 this hyperplane, and  $\xi_i$  the slack variable for sample  $\mathbf{x}_i$ . Two constrains  
 156 ( $\sum_{p=1}^m \gamma_p = 1$  and  $\gamma_p \geq 0$ ) are imposed to ensure that i) the final kernel  
 157  $k(\cdot, \cdot, \boldsymbol{\gamma})$  will not become unbounded due to an arbitrarily large  $\gamma_p$  value; and  
 158 ii)  $k(\cdot, \cdot, \boldsymbol{\gamma})$  will maintain its positive-definiteness.

## 159 2.2. The radius of the minimum enclosing ball (MEB)

160 The MEB is a ball enclosing training samples with the minimum radius.  
 161 Following [25], this ball can be found by solving

$$\min_{R, \mathbf{c}} R^2; \quad s.t. \quad \|\phi(\mathbf{x}_i; \boldsymbol{\gamma}) - \mathbf{c}\|^2 \leq R^2, \ \forall i = 1, \dots, n, \tag{3}$$

163 where  $\mathbf{c}$  and  $R$  denote the centre and the radius. Since  $\phi(\cdot; \boldsymbol{\gamma})$  is usually not  
 164 explicitly known, this problem is solved in its dual form [27]. Let  $\mathbf{K}(\boldsymbol{\gamma})$  be  
 165 the kernel matrices for the final kernel. The dual problem is expressed as

$$R_0^2(\boldsymbol{\gamma}) = \left\{ \max_{\boldsymbol{\beta}} \left[ \boldsymbol{\beta}^\top \text{diag}(\mathbf{K}(\boldsymbol{\gamma})) - \boldsymbol{\beta}^\top \mathbf{K}(\boldsymbol{\gamma}) \boldsymbol{\beta} \right]; \ s.t. \ \boldsymbol{\beta}^\top \mathbf{1} = 1, \ \mathbf{0} \leq \boldsymbol{\beta} \right\}, \tag{4}$$



167 where  $\text{diag}(\mathbf{K}(\boldsymbol{\gamma}))$  denotes a column vector consisting of the diagonal entries  
 168 of  $\mathbf{K}(\boldsymbol{\gamma})$ ,  $\boldsymbol{\beta}$  is an  $n$ -dimensional column vector representing the dual variables,  
 169  $\mathbf{1}$  denotes a column vector with all entries being “1”, and  $R_0(\boldsymbol{\gamma})$  denotes the  
 170 radius of the found MEB. Eq. (4) indicates two things: i) the maximisation  
 171 problem within the curly brackets is the dual problem to solve; and ii) the  
 172 maximum objective function value is equivalent to the squared radius of the  
 173 MEB [25]. Readers are referred to Section 7.1 in [25] or Section 7.3 in the  
 174 SVM tutorial<sup>2</sup> for the detailed derivation of the radius of MEB.

### 175 2.3. Existing Radius-incorporated MKL

176 It is known that the generalisation performance of SVM can be unbiasedly  
 177 estimated by the leave-one-out error (LOO) on a training sample set [28].  
 178 Also, the LOO error is upper bounded by the quantity  $R^2/\rho^2$ , which is known  
 179 as the radius-margin bound (RMB) in the literature [22]. Since the margin  
 180  $\rho$  equals  $1/\|\boldsymbol{\omega}\|$  in SVM, the bound is often expressed as  $R^2\|\boldsymbol{\omega}\|^2$ . To obtain  
 181 a classifier with excellent generalisation performance, a small LOO error is  
 182 desired, which in turn prefers small  $\|\boldsymbol{\omega}\|^2$  and small  $R^2$ . Existing radius-  
 183 incorporated MKL algorithms just implement this idea, in a variety of ways.

184 The algorithm in [18] minimises the following objective function:

$$185 \min_{\boldsymbol{\gamma}, \boldsymbol{\omega}, b, \boldsymbol{\xi}} \frac{1}{2} \|\boldsymbol{\omega}\|^2 + \frac{C}{\sum_{p=1}^m \gamma_p R_p^2} \sum_{i=1}^n \xi_i^2; \quad s.t. \quad y_i(\boldsymbol{\omega}^\top \phi(\mathbf{x}_i; \boldsymbol{\gamma}) + b) \geq 1 - \xi_i, \forall i; \sum_{p=1}^m \gamma_p = 1, \gamma_p \geq 0, \forall p, \quad (5)$$

186 where  $R_p$  is the radius of the MEB in the space  $\mathcal{H}_p$  corresponding to  $k_p(\cdot, \cdot)$ .  
 187 Instead of explicitly computing the radius  $R_0^2(\boldsymbol{\gamma})$  in Eq. (4), this algorithm

---

<sup>2</sup>Available at: <http://research.microsoft.com/pubs/67119/svmtutorial.pdf>

188 computes  $R_p$  for each base kernel and uses its linear combination  $\sum_{p=1}^m \gamma_p R_p^2$   
189 to approximate  $R_0^2(\boldsymbol{\gamma})$ . This approximation avoids solving  $R_0^2(\boldsymbol{\gamma})$  and there-  
190 fore brings computational advantage. In this paper, this approximation  
191 will be adopted into our algorithm. By doing so, we can obtain a radius-  
192 incorporated MKL algorithm that is theoretically more elegant and compu-  
193 tationally more efficient than the one developed in [18].

194 Another algorithm in [19] minimises the following objective function

$$\min_{\boldsymbol{\gamma}, \boldsymbol{\omega}, b, \boldsymbol{\xi}} \frac{1}{2} R_0^2(\boldsymbol{\gamma}) \|\boldsymbol{\omega}\|^2 + C \sum_{i=1}^n \xi_i; \text{ s.t. } y_i(\boldsymbol{\omega}^\top \phi(\mathbf{x}_i; \boldsymbol{\gamma}) + b) \geq 1 - \xi_i, \xi_i \geq 0, \forall i; \gamma_p \geq 0, \forall p. \quad (6)$$

196 As seen,  $R_0^2(\boldsymbol{\gamma})$  is explicitly computed here. To solve this problem, the work  
197 in [19] firstly converts the primal problem to its dual problem as

$$\max_{\boldsymbol{\alpha}} \left\{ \boldsymbol{\alpha}^\top \mathbf{1} - \frac{1}{2R_0^2(\boldsymbol{\gamma})} (\boldsymbol{\alpha} \circ \mathbf{y})^\top \mathbf{K}(\boldsymbol{\gamma}) (\boldsymbol{\alpha} \circ \mathbf{y}) \right\}; \text{ s.t. } \boldsymbol{\alpha}^\top \mathbf{y} = 0; \mathbf{0} \leq \boldsymbol{\alpha} \leq C\mathbf{1}, \quad (7)$$

199 Due to the strong duality [27] of the primal problem, the minimum objective  
200 function value of the primal problem equals the maximum objective function  
201 value of the dual problem. Therefore, the work in [19] defines  $\mathcal{J}(\boldsymbol{\gamma})$  as the  
202 maximum objective function value of Eq. (7) and reformulate Eq. (6) into

$$\min_{\boldsymbol{\gamma}} \mathcal{J}(\boldsymbol{\gamma}); \text{ s.t. } \gamma_p \geq 0, \forall p, \quad (8)$$

204 where  $\mathcal{J}(\boldsymbol{\gamma}) = \max_{\boldsymbol{\alpha}} \left\{ \boldsymbol{\alpha}^\top \mathbf{1} - \frac{1}{2R_0^2(\boldsymbol{\gamma})} (\boldsymbol{\alpha} \circ \mathbf{y})^\top \mathbf{K}(\boldsymbol{\gamma}) (\boldsymbol{\alpha} \circ \mathbf{y}) \right\}$  subject to the  
205 constraints  $\boldsymbol{\alpha}^\top \mathbf{y} = 0$  and  $\mathbf{0} \leq \boldsymbol{\alpha} \leq C\mathbf{1}$ .

206 The work in [19] proposes a tri-level optimisation process to solve the  
207 above problem. Specifically, at each iteration, i) When a new set of kernel  
208 weights  $\boldsymbol{\gamma}$  is obtained,  $R_0^2(\boldsymbol{\gamma})$  will be updated by solving the quadratic pro-  
209 gramming (QP) problem in Eq.(4); ii) The updated  $R_0^2(\boldsymbol{\gamma})$  is combined into

210 Eq.(7) to solve another QP problem to update  $\alpha$  and then a new value of  
211  $\mathcal{J}(\gamma)$  is obtained; iii) After that,  $\gamma$  will be updated according to Eq. (8) and  
212 then go back to Step i). The above procedure is repeated until convergence.  
213 As pointed out in [19], the optimisation problem in Eq. (8) is not convex  
214 with respect to  $\gamma$  and can only converge to a locally optimal solution.

215 Our algorithm will have an optimisation process similar to the above.  
216 However, it does not need to solve the QP problem in Step i) due to adopting  
217 the approximation in [18]. As will be seen, our algorithm can achieve higher  
218 performance in terms of classification and computation than [19]. Last but  
219 no the least, comparing Eqs. (6) and (5) can find that the work in [19] does  
220 not impose the constraint  $\sum_{p=1}^m \gamma_p = 1$  encountered at the end of Section 2.1.  
221 As highlighted in [19], this is an important advantage over [18] because this  
222 constraint is not necessarily an optimal setting. We shall let MKL automati-  
223 cally decide the scale of  $\gamma$  by utilising the radius information. Note that this  
224 desirable property is preserved in our algorithm.

### 225 3. The Proposed Algorithm

226 As indicated in Section 1, we believe that two issues, a non-convex objec-  
227 tive function and the cross-validation-based regularisation parameter tuning,  
228 lead to the unsatisfying prediction performance of existing radius-incorporated  
229 MKL algorithms [18, 19] on AD prediction tasks. Our improvement consists  
230 of two ideas. One is to approximate the radius of MEB with a linear function  
231 of the kernel weights, and the other is to jointly optimise the regularisation  
232 parameter with the kernel weights via MKL.

233 3.1. Approximation to the radius of MEB

234 We firstly show the relationship between the radius of MEB in the feature  
 235 space  $\mathcal{H}(\boldsymbol{\gamma})$  and a linear combination of the  $p$  radiuses in the feature space  
 236  $\mathcal{H}_1, \dots, \mathcal{H}_p$ . It is this relationship who inspires and justifies our idea.

237 This relationship was observed in [18], as stated in Theorem 1. Recall  
 238 that  $\gamma_p$  ( $\gamma_p \geq 0$ ) is the weight for the  $p$ -th base kernel;  $R_0(\boldsymbol{\gamma})$  is the radius of  
 239 MEB defined in Eq.(4); and  $R_p$  is the radius in  $\mathcal{H}_p$ .

240 **Theorem 1.** *It can be proved that  $R_0^2(\boldsymbol{\gamma}) \leq \sum_{p=1}^m \gamma_p R_p^2$ .*

241 For self-containedness, we have provided a more rigorous proof in Appendix.

242 This theorem indicates that  $\sum_{p=1}^m \gamma_p R_p^2$  is an upper bound of  $R_0^2(\boldsymbol{\gamma})$ . Note  
 243 that  $R_p^2$  can be pre-computed once the  $p$  base kernels are predefined, and it  
 244 remains constant. Therefore, the only variables in  $\sum_{p=1}^m \gamma_p R_p^2$  are the weights  
 245  $\boldsymbol{\gamma}$ , and this upper bound is consequently a linear function of  $\boldsymbol{\gamma}$ .

246 Inspired by this observation, we propose to approximate  $R_0^2(\boldsymbol{\gamma})$  with this  
 247 upper bound. It will bring forth the following benefits (as will be shown  
 248 in the following parts): i) Our objective function can be proved to be an  
 249 upper bound of the LOO error. This provides the theoretical justification for  
 250 minimising our objective function to optimise  $\boldsymbol{\gamma}$ ; ii) This approximation is  
 251 pivotal for linking our algorithm to SimpleMKL, which allows it to be readily  
 252 implemented by the off-the-shelf packages; and iii) It makes our optimisation  
 253 problem convex with respect to  $\boldsymbol{\gamma}$ , and this greatly facilities the optimisation.

254 Since this approximation was also used in [18], we highlight the differences  
 255 as follows: Firstly, that work cannot automatically handle the scaling issue  
 256 of  $\boldsymbol{\gamma}$ , as pointed out in [19]. As a result, an additional norm-constraint,  
 257  $\sum_{p=1}^m \gamma_p = 1$ , has to be imposed. However, which norm should be used is an

258 issue itself. Our algorithm is free of this issue; Secondly, our algorithm can  
 259 automatically tune the parameter  $C$ , while cross-validation has to be used in  
 260 [18]; Lastly, our algorithm can achieve superior AD prediction performance  
 261 to the algorithm in [18], as will be experimentally demonstrated.

### 262 3.2. The proposed radius-incorporated MKL algorithm (L2BRMKL)

263 The radius-margin bound is derived based on a hard-margin SVM [22].  
 264 Nevertheless, The classes in AD prediction cannot be well separated in gen-  
 265 eral, and a soft-margin SVM is needed. To make this bound still applicable,  
 266 we use a 2-normed soft-margin SVM<sup>3</sup> because it can be rewritten as a hard-  
 267 margin SVM by slightly modifying its kernel matrix from  $\mathbf{K}$  to  $\mathbf{K} + \mathbf{I}/C$ ,  
 268 where  $\mathbf{I}$  is an identity matrix and  $C$  is the regularisation parameter [22, 21].

269 To incorporate the radius information, we could directly minimise the  
 270 radius-margin bound to optimise the weights  $\gamma$ . Applying the radius-margin  
 271 bound gives the following optimisation problem (We assume that the 2-  
 272 normed soft-margin SVM has been rewritten into a hard-margin one)

$$273 \quad \min_{\gamma} \min_{\omega, b} \frac{1}{2} R_0^2(\gamma) \|\omega\|^2; \quad s.t. \quad y_i(\omega^\top \phi(\mathbf{x}_i; \gamma) + b) \geq 1, \forall i, \quad \gamma_p \geq 0, \forall p, \quad (9)$$

274 However, this optimisation shares the same problem of [19]. That is, a direct  
 275 incorporation of the radius makes the problem non-convex, which is prone  
 276 to being trapped into a local solution.

277 To handle this situation, we propose to approximate  $R_0^2(\gamma)$  with the result  
 278 in Theorem 1. Recall that the kernel matrix is modified from  $\mathbf{K}(\gamma)$  to  $\mathbf{K}(\gamma) +$

---

<sup>3</sup>In a 2-normed soft-margin SVM, the power of the slack variable  $\xi_i$  in the SVM object function is set as 2. The detail can be found in [25] (Chapter 7.2).

279  $\mathbf{I}/C$ . Now we define one more base kernel to account for this modification:  
 280  $\gamma_{m+1}$  is defined as  $1/C$ , and a dummy base kernel matrix is defined as the  
 281 identity matrix  $\mathbf{I}$ . In this way, we can utilise MKL to jointly tune  $C$ .

282 Now, we formally propose the objective function of L2BRMKL as

$$\min_{\gamma} \min_{\omega, b} \frac{1}{2} \left( \sum_{p=1}^{m+1} \gamma_p R_p^2 \right) \|\omega\|^2; \quad s.t. \quad y_i(\omega^\top \phi(\mathbf{x}_i; \gamma) + b) \geq 1, \forall i, \quad \gamma_p \geq 0, \forall p. \quad (10)$$

283

284 Its properties are shown through the following propositions.

285 **Proposition 1.** *The objective function in Eq.(10) is an upper bound of the*  
 286 *radius-margin bound in the form of  $\frac{1}{2}R_0^2(\gamma)\|\omega\|^2$ .*

287 *Proof.* By Theorem 1, we can obtain that  $R_0^2(\gamma) \leq \sum_{p=1}^{m+1} \gamma_p R_p^2$ , and it leads  
 288 to  $\frac{1}{2}R_0^2(\gamma)\|\omega\|^2 \leq \frac{1}{2} \left( \sum_{p=1}^{m+1} \gamma_p R_p^2 \right) \|\omega\|^2$ . This completes the proof.  $\square$

289 Proposition 1 indicates that by minimising this objective function, we can  
 290 restrict the value of the radius-margin bound, which will in turn restrict the  
 291 LOO error, an unbiased estimate of the generalisation error of SVM. This  
 292 provides justification for our objective function.

293 In the following part, we show that our optimisation problem in Eq. (10)  
 294 can be addressed via solving a convex optimisation problem. To this end, we  
 295 rewrite the problem in Eq.(10) as

$$\min_{\gamma} \mathcal{J}(\gamma); \quad s.t. \quad \gamma_p \geq 0, \quad \forall p, \quad (11)$$

296 where

$$\mathcal{J}(\gamma) = \left\{ \min_{\omega, b} \frac{1}{2} \left( \sum_{p=1}^{m+1} \gamma_p R_p^2 \right) \|\omega\|^2; \quad s.t. \quad y_i(\omega^\top \phi(\mathbf{x}_i; \gamma) + b) \geq 1, \forall i \right\}. \quad (12)$$

299 The following Theorem 2 shows that the problem in Eq.(11) can be refor-  
 300 mulated into a form similar to SimpleMKL [17]. The mere but critical dif-  
 301 ference is that a radius-weighted norm-constraint is used, compared with an  
 302 unweighted version in[17]. The significance of Theorem 2 lies at that i) It  
 303 reveals the connection of our algorithm with traditional MKL without using  
 304 the radius information; ii) It shows that our radius-incorporated MKL al-  
 305 gorithm, which appears to be sophisticated, can essentially be reduced to a  
 306 slightly changed traditional MKL. iii) It suggests that our algorithm can be  
 307 efficiently solved by the existing MKL software packages.

308 To prove Theorem 2, we first give Proposition 2 for the optimisation  
 309 problem in Eq.(11). Its proof can be found in our previous work [20].

310 **Proposition 2.** *The objective function value  $\mathcal{J}(\boldsymbol{\gamma})$  remains unchanged when  
 311  $\boldsymbol{\gamma}$  is scaled to  $\tau\boldsymbol{\gamma}$ , where the scale factor  $\tau$  is any positive scalar. Also, the  
 312 SVM decision function of the resulting MKL algorithm is not affected by  $\tau$ .*

313 With Proposition 2, Theorem 2 shows that solving the optimisation in  
 314 Eq.(11) can be converted to solving a related but simpler optimisation.

315 **Theorem 2.** *The optimal solution of the optimisation problem in Eq.(11),  
 316 denoted as  $\boldsymbol{\gamma}^*$ , can be written as  $\boldsymbol{\gamma}^* = \left(\sum_{p=1}^{m+1} \gamma_p^* R_p^2\right) \boldsymbol{\eta}^*$ , where  $\boldsymbol{\eta}^*$  is the  
 317 optimal solution of the following optimisation problem*

$$318 \quad \min_{\boldsymbol{\eta}} \mathcal{J}(\boldsymbol{\eta}); \quad s.t. \quad \sum_{p=1}^{m+1} \eta_p R_p^2 = 1, \eta_p \geq 0, \quad \forall p. \quad (13)$$

319 where

$$320 \quad \mathcal{J}(\boldsymbol{\eta}) = \left\{ \min_{\tilde{\boldsymbol{\omega}}, b} \frac{1}{2} \|\tilde{\boldsymbol{\omega}}\|^2; \quad s.t. \quad y_i (\tilde{\boldsymbol{\omega}}^\top \phi(\mathbf{x}_i; \boldsymbol{\eta}) + b) \geq 1, \forall i \right\}. \quad (14)$$

321 Also, for the SVM decision function of the resulting MKL algorithm, denoted  
 322 by  $f(\mathbf{x})$ , it can be proved that  $f(\mathbf{x}) = \tilde{\boldsymbol{\omega}}^\top \phi(\mathbf{x}; \boldsymbol{\eta}) + b = \boldsymbol{\omega}^\top \phi(\mathbf{x}; \boldsymbol{\gamma}) + b$ .

323 *Proof.* Defining  $\tilde{\boldsymbol{\omega}} := \sqrt{(\sum_{p=1}^{m+1} \gamma_p R_p^2)} \boldsymbol{\omega}$ , Eq.(12) is rewritten as

$$324 \quad \mathcal{J}(\boldsymbol{\gamma}) = \left\{ \min_{\tilde{\boldsymbol{\omega}}, b} \frac{1}{2} \|\tilde{\boldsymbol{\omega}}\|^2; \text{ s.t. } y_i \left( \tilde{\boldsymbol{\omega}}^\top \phi \left( \mathbf{x}_i; \frac{\boldsymbol{\gamma}}{\sum_{p=1}^{m+1} \gamma_p R_p^2} \right) + b \right) \geq 1, \forall i. \right\} \quad (15)$$

325 Letting  $\tau = \frac{1}{\sum_{p=1}^{m+1} \gamma_p R_p^2}$  and  $\boldsymbol{\eta} = \tau \boldsymbol{\gamma}$  (that is,  $\eta_p = \tau \gamma_p, \forall p$ ), we obtain that

$$326 \quad \sum_{p=1}^{m+1} \eta_p R_p^2 = \sum_{p=1}^{m+1} \tau \gamma_p R_p^2 = \tau \sum_{p=1}^{m+1} \gamma_p R_p^2 = 1. \quad (16)$$

327 The last equality is a direct result of the definition of  $\tau$ . Also, applying

328 Proposition 2, we can obtain  $\mathcal{J}(\boldsymbol{\eta}) = \mathcal{J}(\tau \boldsymbol{\gamma}) = \mathcal{J}(\boldsymbol{\gamma})$ . Hence, Eq.(15) can

329 be rewritten with respect to  $\boldsymbol{\eta}$  as

$$330 \quad \mathcal{J}(\boldsymbol{\eta}) = \left\{ \min_{\tilde{\boldsymbol{\omega}}, b} \frac{1}{2} \|\tilde{\boldsymbol{\omega}}\|^2; \text{ s.t. } y_i (\tilde{\boldsymbol{\omega}}^\top \phi(\mathbf{x}_i; \boldsymbol{\eta}) + b) \geq 1, \forall i; \sum_{p=1}^{m+1} \eta_p R_p^2 = 1, \eta_p \geq 0, \forall p. \right\}, \quad (17)$$

331 Because  $\eta_p$  is not a variable of this problem, the constraint on  $\eta_p$  can be

332 moved out of the curly brackets and this gives the result in Eq. (14). Then,

333 combing the constraints on  $\eta_p$  with  $\min_{\boldsymbol{\eta}} \mathcal{J}(\boldsymbol{\eta})$  leads to the optimisation

334 problem in Eq.(13) exactly. This proves the first part of this theorem.

335 We now prove the second part on SVM decision function. Note that

336  $\phi(\mathbf{x}; \boldsymbol{\eta})$  can be written as  $[\sqrt{\eta_1} \phi_1(\mathbf{x})^\top, \dots, \sqrt{\eta_{m+1}} \phi_{m+1}(\mathbf{x})^\top]^\top$ . We partition

337  $\tilde{\boldsymbol{\omega}}$  in a similar manner as  $\tilde{\boldsymbol{\omega}} = [\tilde{\boldsymbol{\omega}}_1^\top, \dots, \tilde{\boldsymbol{\omega}}_{m+1}^\top]^\top$ , and then obtain that

$$338 \quad \begin{aligned} f(\mathbf{x}) &= \tilde{\boldsymbol{\omega}}^\top \phi(\mathbf{x}; \boldsymbol{\eta}) + b = \sum_{p=1}^{m+1} \tilde{\boldsymbol{\omega}}_p^\top \sqrt{\eta_p} \phi_p(\mathbf{x}) + b = \sum_{p=1}^{m+1} \tilde{\boldsymbol{\omega}}_p^\top \sqrt{\tau \gamma_p} \phi_p(\mathbf{x}) + b \\ &= \sum_{p=1}^{m+1} \boldsymbol{\omega}_p^\top \sqrt{\gamma_p} \phi_p(\mathbf{x}) + b = \boldsymbol{\omega}^\top \phi(\mathbf{x}; \boldsymbol{\gamma}) + b. \end{aligned} \quad (18)$$

339 In the last two steps, we use the following facts: i)  $\tilde{\boldsymbol{\omega}} := \sqrt{(\sum_{p=1}^{m+1} \gamma_p R_p^2)} \boldsymbol{\omega} =$

340  $\boldsymbol{\omega} / \sqrt{\tau}$ ; and ii)  $\phi(\mathbf{x}; \boldsymbol{\gamma}) = [\sqrt{\gamma_1} \phi_1(\mathbf{x})^\top, \dots, \sqrt{\gamma_{m+1}} \phi_{m+1}(\mathbf{x})^\top]^\top$  and the parti-

341 tion that  $\boldsymbol{\omega} = [\boldsymbol{\omega}_1^\top, \dots, \boldsymbol{\omega}_{m+1}^\top]^\top$ . This completes the proof.  $\square$



342 As shown by Theorem 2, the solution of Eq.(11) can be obtained by solv-  
 343 ing Eq.(13). In the following part, we prove that the optimisation problem in  
 344 Eq.(13) can be reformulated as a convex one. This will justify our claim that  
 345 our algorithm can be addressed by solving a convex optimisation problem.

346 **Proposition 3.** *The optimisation problem in Eq.(13) is equivalent to*

$$\begin{aligned}
 & \min_{\boldsymbol{\eta}} \min_{\widehat{\boldsymbol{\omega}}, b} \frac{1}{2} \sum_{p=1}^{m+1} \frac{\|\widehat{\boldsymbol{\omega}}_p\|^2}{\eta_p} \\
 & \text{s.t. } y_i \left( \sum_{p=1}^{m+1} \widehat{\boldsymbol{\omega}}_p^\top \phi_p(\mathbf{x}_i) + b \right) \geq 1, \forall i; \quad \sum_{p=1}^{m+1} \eta_p R_p^2 = 1, \eta_p \geq 0, \forall p,
 \end{aligned} \tag{19}$$

348 *which is jointly convex with respect to  $\boldsymbol{\eta}$ ,  $\widehat{\boldsymbol{\omega}}$  and  $b$ .*

349 *Proof.* Let us define  $\widehat{\boldsymbol{\omega}}_p := \sqrt{\eta_p} \widetilde{\boldsymbol{\omega}}_p$ . By substituting  $\widehat{\boldsymbol{\omega}}_p$  into Eq. (14) and  
 350 replacing  $\mathcal{J}(\boldsymbol{\eta})$  in Eq.(13) with the result in Eq. (14) after substitution, we  
 351 obtain the optimisation problem in Eq.(19). Its objective function is a ratio  
 352 of a quadratic function ( $\|\widehat{\boldsymbol{\omega}}_p\|^2$ ) to a linear function ( $\eta_p$ ). According to [27]  
 353 (Chapter 3.2.6, example 3.18 on page 89), this function is convex. Also,  
 354 because all the constraints are linear function of  $\boldsymbol{\eta}$ ,  $\widehat{\boldsymbol{\omega}}$  and  $b$ , the feasible  
 355 domain of this optimisation is convex. Therefore, the problem in Eq.(19) is  
 356 convex with respect to its variables. This completes the proof.  $\square$

357 Interestingly, we find that Eq.(19) has a form similar to the one in [17],  
 358 with the only difference that the constraint  $\sum_{p=1}^{m+1} \eta_p R_p^2 = 1$  is used in our  
 359 algorithm while  $\sum_{p=1}^m \eta_p = 1$  is used in [17]. The problem in Eq.(19) can be  
 360 solved by any MKL packages sharing the same routine: updating the struc-  
 361 tural parameters of SVM,  $\boldsymbol{\alpha}$ , and the weights of base kernels,  $\boldsymbol{\eta}$ , alternately.

362 Specifically,  $\boldsymbol{\alpha}$  is updated with the current  $\boldsymbol{\eta}$  by solving

$$363 \quad \max_{\boldsymbol{\alpha}} \left\{ \mathbf{1}^\top \boldsymbol{\alpha} - \frac{1}{2} (\boldsymbol{\alpha} \circ \mathbf{y})^\top \mathbf{K}(\boldsymbol{\eta}) (\boldsymbol{\alpha} \circ \mathbf{y}) \right\}; \quad s.t. \quad \boldsymbol{\alpha}^\top \mathbf{y} = 0, \quad \boldsymbol{\alpha} \geq \mathbf{0}, \quad (20)$$

364 which is the dual problem of Eq.(14). Then, the weights  $\boldsymbol{\eta}$  are updated  
365 with the obtained  $\boldsymbol{\alpha}$  by using the reduced gradient descent method [17].

366 Specifically, the cost function for updating  $\boldsymbol{\eta}$  is

$$367 \quad \min_{\boldsymbol{\eta}} \mathbf{1}^\top \boldsymbol{\alpha}_0 - \frac{1}{2} (\boldsymbol{\alpha}_0 \circ \mathbf{y})^\top \mathbf{K}(\boldsymbol{\eta}) (\boldsymbol{\alpha}_0 \circ \mathbf{y}) \quad s.t. \quad \sum_{p=1}^{m+1} \eta_p R_p^2 = 1, \quad \eta_p \geq 0, \quad \forall p, \quad (21)$$

368 where  $\boldsymbol{\alpha}_0$  is obtained in the last iteration with fixed  $\boldsymbol{\eta}$ . Note that Eq.(21)  
369 is a constrained optimisation problem w.r.t  $\boldsymbol{\eta}$ . The positivity and equality  
370 constraints have to be maintained during the update of  $\boldsymbol{\eta}$ . Such problems  
371 can be effectively solved via the reduced gradient descent method [17]. This  
372 procedure repeats until convergence. Our algorithm is listed in Algorithm 1.

---

**Algorithm 1** The proposed L2BRMKL

---

- 1: Initialise  $\boldsymbol{\eta}^0$ ,  $\mathbf{K}_{m+1} = \mathbf{I}$  and  $\eta_{m+1} = 1/C_0$ .
  - 2: Calculate  $R_p^2$  ( $p = 1, \dots, m+1$ ) for each base kernel by following Eq.(3).
  - 3:  $i \leftarrow 0$
  - 4: **repeat**
  - 5:     Obtain  $\boldsymbol{\alpha}^{i+1}$  by solving Eq.(20) with  $\boldsymbol{\eta}^i$ .
  - 6:     Update  $\boldsymbol{\eta}^{i+1}$  by solving Eq.(21) via the reduced gradient descent method [17] with  $\boldsymbol{\alpha}^{i+1}$ .
  - 7:      $i \leftarrow i + 1$
  - 8: **until**  $\max \{ |\eta_1^{i+1} - \eta_1^i|, \dots, |\eta_{m+1}^{i+1} - \eta_{m+1}^i| \} < 1e - 4$
- 

373

374 Applying Proposition 2 (or Theorem 2), we know that the SVM decision  
375 function produced by the problems in Eqs.(11) and (13) are the same. Also,

376 according to Proposition 3, we can solve Eq.(13) by solving the equivalent  
 377 problem in Eq.(19). Therefore, after we obtain the optimal  $\boldsymbol{\alpha}^*$ ,  $b^*$  and  $\boldsymbol{\eta}^*$  by  
 378 solving Eq.(19), we can directly write the SVM decision function as

$$379 \quad f(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y_i \sum_{p=1}^m \eta_p^* k_p(\mathbf{x}_i, \mathbf{x}) + b^*, \quad (22)$$

380 and it will be used to classify test samples.

381 We close this subsection by highlighting the differences between the pro-  
 382 posed L2BRMKL and our previous  $\ell_2\text{tr}\mathbf{S}_t\text{MKL}$  in [20], where the trace of a  
 383 total scatter matrix in the feature space  $\mathcal{H}(\boldsymbol{\gamma})$  is used to substitute  $R_0^2(\boldsymbol{\gamma})$ .  
 384 Compared with  $\ell_2\text{tr}\mathbf{S}_t\text{MKL}$ , L2BRMKL has an objective function that can  
 385 be proved as an upper bound related to the generalisation performance of  
 386 the SVM classifier, as in Proposition 1. However,  $\ell_2\text{tr}\mathbf{S}_t\text{MKL}$  does not have  
 387 this merit because the trace of the total scatter matrix is only a lower (rather  
 388 than upper) bound of  $R_0^2(\boldsymbol{\gamma})$  [29]. This difference is important in that it pro-  
 389 vides better theoretical justification for our algorithm. And we do observe  
 390 the improvement brought by such a difference in the experimental study.

## 391 4. Discussion on the Improvement

392 In this section, we discuss the potential aspects that contribute to the  
 393 improvement achieved by our algorithm. Also, to better demonstrate how  
 394 these aspects contribute, we develop two additional MKL algorithms in which  
 395 only one aspect is improved while the other is kept unchanged.

### 396 4.1. Approximating the radius

397 To show the help of this aspect, we modify the objective function of the  
 398 algorithm in [19], which directly incorporates the radius information by using

399  $R_0^2(\boldsymbol{\gamma})$ . Maintaining all the other settings in [19], we only replace  $R_0^2(\boldsymbol{\gamma})$  with  
 400  $\sum_{p=1}^m \gamma_p R_p^2$ . Specifically, its objective function now becomes

$$\begin{aligned}
 & \min_{\boldsymbol{\gamma}, \boldsymbol{\omega}, b, \boldsymbol{\xi}} \frac{1}{2} \left( \sum_{p=1}^m \gamma_p R_p^2 \right) \|\boldsymbol{\omega}\|^2 + C \sum_{i=1}^n \xi_i \\
 & \text{s.t. } y_i(\boldsymbol{\omega}^\top \phi(\mathbf{x}_i; \boldsymbol{\gamma}) + b) \geq 1 - \xi_i, \xi_i \geq 0, \forall i, \sum_{p=1}^m \gamma_p = 1, \gamma_p \geq 0, \forall p.
 \end{aligned}
 \tag{23}$$

402 Following the work in [19], the parameter  $C$  is chosen by cross-validation.  
 403 Since this objective function is based on the SVM classifier with 1-normed  
 404 soft margin, we call the resulting algorithm L1BRMKL+C. By comparing  
 405 this algorithm with that in [19], we can see whether the approximation to  
 406 the radius contributes to the improvement on the performance of MKL.

#### 407 4.2. Automatically tuning the regularisation parameter $C$

408 We design another algorithm, termed as L2BRMKL+C, to investigate the  
 409 benefit of automatically tuning  $C$  in MKL algorithms. L2BRMKL+C shares  
 410 the same optimisation problem with the proposed L2BRMKL. They only  
 411 differ in that L2BRMKL+C determines  $C$  by employing multi-fold cross-  
 412 validation still. By comparing L2BRMKL with L2BRMKL+C, we can see  
 413 whether automatically tuning  $C$  really helps.

### 414 5. Experimental Result

415 This experiment aims to evaluate the proposed MKL algorithm, L2BRMKL,  
 416 with respect to classification accuracy and computational efficiency. It is  
 417 compared to a set of state-of-the-art MKL algorithms, including i) the com-  
 418 monly used margin-only algorithm SimpleMKL [17]; ii) three existing radius-  
 419 incorporated algorithms RMKL [18], MBMKL [19], and  $\ell_2\text{tr}\mathbf{S}_t$ MKL [20]; iii)

420 recently developed non-sparse MKL algorithm NSMKL [24] which constrain  
421 the kernel weights with different norms; and iv) unweighted MKL algorithm  
422 UWMKL that simply uses the average of all base kernels.

### 423 5.1. Data sets and experimental settings

424 We firstly use the 11 UCI machine learning benchmark data sets, which  
425 have been widely used to evaluate MKL algorithms [19, 18, 24]. Their names  
426 are listed in Table 2, and the data sets can be downloaded from the Internet<sup>4</sup>.

427 Every feature in these data sets is normalised to have zero mean and  
428 unit variance. To accumulate statistic, 30 training and test splits are cre-  
429 ated for each data set. For each split, 20% of samples in the data set are  
430 randomly selected as training data and the rest 80% is used for test. To pre-  
431 define base kernels, we adopt four types of kernel functions, including Gaus-  
432 sian kernel  $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/\sigma)$ , Laplacian kernel  $k(\mathbf{x}_i, \mathbf{x}_j) =$   
433  $\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|/\sqrt{\sigma})$ , Inverse square distance kernel  $k(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{\|\mathbf{x}_i - \mathbf{x}_j\|^2/\sigma+1}$ ,  
434 and Inverse distance kernel  $k(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{\|\mathbf{x}_i - \mathbf{x}_j\|/\sqrt{\sigma+1}}$ , where  $\sigma$  denotes the ker-  
435 nel parameter. Let  $\sigma_0$  stand for the average value of the pairwise Euclidean  
436 distance between samples in a training set. In this experiment,  $\sigma$  is set as  $2^t\sigma_0$   
437 with  $t = -2, -1, 0, 1, 2$ , respectively, and employed for each kernel. In this  
438 way, we generate 20 ( $4 \times 5$ ) base kernels and use them in all the algorithms  
439 conducted on the 11 UCI data sets.

440 For each data set, an algorithm is trained and test on the 30 training/test  
441 splits, and the average classification accuracy and standard deviation are  
442 reported. To conduct a rigorous comparison, the *paired Student's t-test* is

---

<sup>4</sup><http://archive.ics.uci.edu/ml/datasets.html>

443 performed. The  $p$ -value of this test represents the probability that the two  
444 sets of results in comparison come from the distributions with an equal mean.  
445 A  $p$ -value of 0.05 is considered statistically significant and used here.

446 In addition to the UCI data sets, an AD data set from the Alzheimer  
447 Disease Neuroimaging Initiative (ADNI) database is used. Each sample has  
448 229 features from four data sources, including three CSF (cerebrospinal fluid)  
449 biomarkers, 63 left hippocampal shape features, 63 right hippocampal shape  
450 features and 100 regional grey matter volumes. As shown in Table 1, three  
451 tasks are defined to differentiate different clinic groups, including MCI versus  
452 NC (normal control), PMCI (converters) versus NC, and PMCI versus SMCI  
453 (non-converters). In general, these clinic groups are largely overlapped with  
454 each other, making the prediction tasks challenging. MCI consists of two  
455 subgroups PMCI and SMCI. PMCI is more AD-like, while SMCI is more  
456 NC-like. Therefore, there is an increasing degree of difficulty to differenti-  
457 ate PMCI from NC, MCI from NC, and PMCI from SMCI. Especially, the  
458 last task, which tells the MCI converters (PMCI) from the non-converters  
459 (SMCI), is very challenging and also of great importance in AD prediction.

460 To predefine the base kernels, this experiment applies the above four  
461 types of kernel functions to each data source. For each source, five different  
462 kernel parameter  $\sigma$  values are set for each type of kernel in the same way  
463 as the UCI data sets. By doing so, 80 ( $4 \times 4 \times 5$ ) base kernels are created  
464 in total, with 20 base kernels for each data source. As seen from Table 1,  
465 the number of samples in these prediction tasks is generally small. To make  
466 a good use of these samples, we employ the LOO strategy to evaluate each  
467 MKL algorithm. In this strategy, each sample is used as the test sample

Table 1: Summary of the data sets used in the AD prediction experiments.

Data set	Instances		Features			
	# Positive	# Negative	# CSF biomarkers	hippocampal shape		# regional gray matter volumes
				# left	# right	
PMCI vs. NC	50	70	3	63	63	100
MCI vs. NC	121	70	3	63	63	100
PMCI vs. SMCI	50	71	3	63	63	100

468 in turn to form a classification session. The classification results of all the  
 469 sessions are pooled to obtain classification accuracy.

470 The proposed L2BRMKL algorithm can automatically tune  $C$ . For the  
 471 other algorithms,  $C$  has to be chosen by cross-validation. In this experiment,  
 472 four-fold cross-validation is applied to a large range  $[2^{-5}, 2^{-3}, \dots, 2^{15}]$  to  
 473 choose  $C$ . The experiment is conducted on a high-performance cluster, where  
 474 each node has eight cores with 2.3GHz CPU and 2GB memory.

475 *5.2. Results on UCI data sets*

476 The classification results are listed in Table 2. For each data set, the  
 477 highest accuracy and those whose differences from the highest one are not  
 478 statistically significant are shown in bold. From this table, we can observe  
 479 i) The proposed L2BRMKL and the existing radius-incorporated MKL algo-  
 480 rithms [18, 19, 20] (with average accuracy of 77.5%, 75.9%, 77.0% and 76.1%)  
 481 achieve overall better classification than the margin-based ones [17, 24] (with  
 482 average accuracy of 75.8% and 72.8% or 72.7%). This demonstrates the ef-  
 483 fectiveness of incorporating the radius information; ii) Among the radius-  
 484 incorporated MKL algorithms, L2BRMKL attains the highest average ac-  
 485 curacy 77.5% and wins on nine of the 11 data sets. This result initially  
 486 validates its advantage and provides a basis to investigate its performance

487 on AD prediction tasks in further.

### 488 5.3. Results on the tasks of AD prediction

489 In addition to classification accuracy, we adopt another criterion widely  
490 used in medical applications, i.e., Matthews Correlation Coefficient (MCC),  
491 to evaluate the proposed algorithm. MCC is defined as

$$492 \quad \text{MCC} = \frac{TP * TN - FP * FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}, \quad (24)$$

493 where  $TP$ ,  $TN$ ,  $FP$  and  $FN$  represent the number of true positives, the  
494 number of true negatives, the number of false positives and the number of  
495 false negatives, respectively. As seen from the definition, MCC takes both  
496 “sensitivity” and “specificity” of the classification into account, and is a  
497 balanced measure of classification performance. The classification accuracy  
498 and MCC results are reported in Tables 3 and 4, respectively. As previous,  
499 the highest values are highlighted in bold for each task. Note that no standard  
500 deviation is reported here, because there is only one classification accuracy  
501 for each task when the LOO strategy is used. Also, due to this fact, the  
502 statistical test becomes inapplicable and no result is reported either.

503 From Table 3, we can see that L2BRMKL achieves the overall best clas-  
504 sification performance in the three tasks. In the task of PMCI vs. NC, it  
505 obtains the equally best performance as  $\ell_2\text{tr}\mathbf{S}_t\text{MKL}$ , RMKL and MBMKL.  
506 In the other two tasks, it outperforms all the other algorithms. Also, the  
507 advantage of L2BRMKL becomes more pronounced with the increasing de-  
508 gree of difficulty of the tasks. As seen, the most significant improvement is  
509 achieved on the task of PMCI vs. SMCI, where 3.3% improvement (75.2%  
510 vs. 71.9%) is gained over the second best one, SimpleMKL. In addition, for



Table 2: Comparison of classification accuracy (in percentage) obtained by different MKL algorithms on UCI data sets and the statistical test result. Boldface indicates the highest accuracy and those whose differences from the highest accuracy are not statistically significant (evaluated by paired Student’s  $t$ -test with  $p$ -value  $\geq 0.05$ ). The three numbers in each cell represent the average classification accuracy, standard deviation and the  $p$ -value.

Data set	L2BRMKL	$\ell_2$ trStMKL	SimpleMKL	RMKL	MBMKL	Non-SparseMKL [24]		UWMKL
	<b>Proposed</b>	[20]	[17]	[18]	[19]	$p = 2$	$p = 3$	
Coloncancer	<b>71.6</b>	<b>71.5</b>	67.9	68.5	68.3	66.7	66.0	68.1
	$\pm 6.6$	$\pm 6.8$	$\pm 6.3$	$\pm 7.0$	$\pm 7.1$	$\pm 5.6$	$\pm 4.8$	$\pm 5.2$
	<b>1.00</b>	<b>0.77</b>	0.00	0.00	0.00	0.00	0.00	0.00
Fourclass	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>	<b>99.9</b>	<b>99.8</b>	99.9	99.0
	$\pm 0.0$	$\pm 0.1$	$\pm 0.0$	$\pm 0.1$	$\pm 0.0$	$\pm 0.1$	$\pm 0.2$	$\pm 0.5$
	<b>1.00</b>	<b>0.16</b>	<b>1.00</b>	<b>0.16</b>	<b>1.00</b>	<b>0.06</b>	0.02	0.00
Germannum	71.6	71.5	71.1	71.5	70.9	<b>72.9</b>	<b>72.7</b>	72.2
	$\pm 1.0$	$\pm 1.2$	$\pm 1.3$	$\pm 1.5$	$\pm 1.2$	$\pm 1.5$	$\pm 1.5$	$\pm 1.7$
	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>	<b>0.31</b>	0.02
Heart	<b>82.9</b>	81.5	<b>82.2</b>	80.9	<b>82.7</b>	81.4	81.1	79.5
	$\pm 1.6$	$\pm 2.4$	$\pm 2.5$	$\pm 2.6$	$\pm 1.6$	$\pm 3.1$	$\pm 2.9$	$\pm 3.6$
	<b>1.00</b>	0.00	<b>0.15</b>	0.00	<b>0.41</b>	0.01	0.00	0.00
Ionosphere	<b>66.2</b>	64.7	65.1	<b>65.3</b>	<b>66.4</b>	65.0	65.0	65.1
	$\pm 3.0$	$\pm 1.1$	$\pm 1.2$	$\pm 1.5$	$\pm 3.3$	$\pm 1.2$	$\pm 1.2$	$\pm 1.2$
	<b>0.52</b>	0.00	0.02	<b>0.08</b>	<b>1.00</b>	0.02	0.02	0.03
Liver	<b>60.6</b>	<b>59.9</b>	59.6	59.3	<b>60.4</b>	<b>61.2</b>	<b>61.5</b>	<b>61.7</b>
	$\pm 2.3$	$\pm 2.7$	$\pm 2.5$	$\pm 2.5$	$\pm 2.6$	$\pm 3.6$	$\pm 3.9$	$\pm 4.4$
	<b>0.22</b>	<b>0.05</b>	0.01	0.00	<b>0.16</b>	<b>0.30</b>	<b>0.69</b>	<b>1.00</b>
Musk1	<b>83.1</b>	76.7	78.9	80.4	82.2	53.8	53.3	51.2
	$\pm 3.1$	$\pm 8.9$	$\pm 4.5$	$\pm 4.2$	$\pm 3.7$	$\pm 5.9$	$\pm 6.9$	$\pm 6.4$
	<b>1.00</b>	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Sonar	<b>75.9</b>	<b>75.7</b>	73.2	<b>74.9</b>	74.6	<b>74.2</b>	<b>74.0</b>	<b>73.9</b>
	$\pm 3.2$	$\pm 3.1$	$\pm 4.5$	$\pm 4.5$	$\pm 2.9$	$\pm 4.3$	$\pm 5.6$	$\pm 5.7$
	<b>1.00</b>	<b>0.65</b>	0.00	<b>0.11</b>	0.00	<b>0.06</b>	<b>0.12</b>	<b>0.10</b>
Splice	68.7	63.7	64.7	62.1	<b>70.4</b>	56.2	55.8	55.1
	$\pm 5.2$	$\pm 3.9$	$\pm 5.4$	$\pm 4.4$	$\pm 4.5$	$\pm 4.4$	$\pm 4.6$	$\pm 4.5$
	0.00	0.00	0.00	0.00	<b>1.00</b>	0.00	0.00	0.00
Wdbc	<b>95.9</b>	<b>95.8</b>	95.4	<b>95.8</b>	<b>95.6</b>	94.1	94.2	94.4
	$\pm 0.9$	$\pm 1.0$	$\pm 1.1$	$\pm 1.0$	$\pm 1.1$	$\pm 2.1$	$\pm 2.4$	$\pm 2.3$
	<b>1.00</b>	<b>0.52</b>	0.03	<b>0.58</b>	<b>0.15</b>	0.00	0.00	0.00
Wpbc	<b>76.0</b>	<b>75.8</b>	<b>76.0</b>	<b>76.1</b>	<b>76.1</b>	<b>75.6</b>	<b>75.7</b>	<b>75.3</b>
	$\pm 0.3$	$\pm 1.3$	$\pm 0.5$	$\pm 0.6$	$\pm 0.2$	$\pm 1.8$	$\pm 1.4$	$\pm 2.8$
	<b>0.23</b>	<b>0.14</b>	<b>0.21</b>	<b>1.00</b>	<b>0.60</b>	<b>0.12</b>	<b>0.16</b>	<b>0.15</b>
Average	<b>77.5</b>	76.1	75.8	75.9	77.0	72.8	72.7	72.3
Win	<b>9</b>	6	3	5	7	5	4	3

Table 3: Classification accuracy (in percentage) of different MKL algorithms.

Data set	L2BRMKL	$\ell_2\text{trStMKL}$	SimpleMKL	RMKL	MBMKL	Non-SparseMKL [24]		UWMKL
	Proposed	[20]	[17]	[18]	[19]	$p = 2$	$p = 3$	
PMCI vs. NC	<b>88.3</b>	<b>88.3</b>	87.5	<b>88.3</b>	<b>88.3</b>	85.8	85.0	85.8
MCI vs. NC	<b>75.4</b>	74.4	72.3	69.6	72.3	74.8	73.8	73.3
PMCI vs. SMCI	<b>75.2</b>	71.1	71.9	67.8	69.4	70.3	67.8	67.8
Average	<b>79.6</b>	77.9	77.2	75.2	76.7	77.0	75.5	75.6
Win	<b>3</b>	1	0	1	1	0	0	0

Table 4: Comparison of MCC (in percentage) of different MKL algorithms.

Data set	L2BRMKL	$\ell_2\text{trStMKL}$	SimpleMKL	RMKL	MBMKL	Non-SparseMKL [24]		UWMKL
	Proposed	[20]	[17]	[18]	[19]	$p = 2$	$p = 3$	
PMCI vs. NC	<b>76.2</b>	<b>76.2</b>	74.9	76.0	<b>76.2</b>	71.2	69.7	71.0
MCI vs. NC	<b>45.9</b>	43.7	38.0	30.9	38.8	44.7	43.6	42.3
PMCI vs. SMCI	<b>48.2</b>	39.9	40.9	32.1	35.5	38.3	33.3	33.0
Average	<b>56.8</b>	53.3	51.3	46.3	50.2	51.4	48.9	48.8
Win	<b>3</b>	1	0	0	1	0	0	0

511 the easiest task of PMCI vs. NC, all the four radius-incorporated methods  
512 (L2BRMKL, RMKL, MBMKL, and  $\ell_2\text{trS}_t\text{MKL}$ ) perform better than the  
513 margin-only methods, indicating the advantage of incorporating radius in-  
514 formation again. With the introduction of SMCI, the second and third tasks  
515 become more difficult. As observed, the performance of  $\ell_2\text{trS}_t\text{MKL}$ , RMKL  
516 and MBMKL decreases significantly from 88.3% to 71.1%, 88.3% to 67.8%,  
517 and 88.3% to 69.4%, respectively. Although the performance of L2BRMKL  
518 also decreases due to the increased difficulty, its decrease is the smallest one.  
519 These results demonstrate the superiority of the proposed L2BRMKL on  
520 these AD prediction tasks. Also, this situation is further confirmed by the  
521 MCC values reported in Table 4, where the proposed L2BRMKL consistently  
522 shows the highest MCC values on the three AD prediction tasks.

523 To check whether the approximation to the radius (discussed in Sec-  
 524 tion 4.1) contributes to the above improvement, we compare MBMKL in [19]  
 525 and the L1BRMKL+C in Section 4. Recall that they only differ in the way  
 526 to estimate the radius. The results are reported in Table 5. The performance  
 of L2BRMKL is also quoted for reference. From this table, we can see that

Table 5: Classification accuracy (in %) of L2BRMKL, L1BRMKL+C and MBMKL [19].

Data set	L2BRMKL Proposed	L1BRMKL+C	MBMKL [19]
PMCI vs. NC	88.3	<b>90.8</b>	88.3
MCI vs. NC	<b>75.4</b>	71.2	72.3
PMCI vs. SMCI	<b>75.2</b>	73.6	69.4
Average	<b>79.6</b>	78.5	76.7

527

528 L1BRMKL+C shows an overall better performance than MBMKL (78.5%  
 529 vs. 76.7% by average), and achieves a clear improvement by 4.2% (73.6% vs.  
 530 69.4%) at the most difficult task of PMCI vs. SMCI. We conjecture that the  
 531 inferiority of MBMKL to L1BRMKL+C is due to the issue of numerical in-  
 532 stability. Recall that MBMKL has a tri-level optimisation process, in which  
 533 the radius of MEB is updated by solving a QP problem at each iteration.  
 534 In this case, any numerical error occurring in solving this QP problem could  
 535 adversely affect and entangle with the optimisation of structural parameter  
 536  $\alpha$  and in turn the kernel weights  $\gamma$ . And such numerical errors could ac-  
 537 cumulate with the increasing number of iterations. A rigorous theoretical  
 538 analysis of this numerical instability issue will be a good topic in our future  
 539 work. At last, note that despite its advantage over MBMKL, L1BRMKL+C  
 540 still performs overall worse than the proposed L2BRMKL.

541 To check the contribution of automatic tuning of the regularisation pa-

parameter  $C$ , we compare L2BRMKL with the L2BRMKL+C in Section 4.2. The only difference between them is how  $C$  is tuned. As shown in Table 6,

Table 6: Classification accuracy (in percentage) of L2BRMKL and L2BRMKL+C.

Data set	L2BRMKL	L2BRMKL+C
	Proposed	
PMCI vs. NC	<b>88.3</b>	87.5
MCI vs. NC	<b>75.4</b>	71.2
PMCI vs. SMCI	<b>75.2</b>	68.6
Average	<b>79.6</b>	75.8

543

L2BRMKL consistently outperforms L2BRMKL+C. Also, the more difficult the task is, the more the improvement is attained. These results demonstrate the benefit of automatically tuning  $C$  on these classification tasks. The inferior performance of L2BRMKL+C is due to the following two factors: i) the estimate of  $C$  obtained by cross-validation becomes unreliable when the number of training samples is small; and ii) cross-validation can only examine a limited number of possible  $C$  values. These could result in a less appropriate  $C$  value and adversely affect the performance of L2BRMKL+C.

#### 552 5.4. Computational efficiency

553 In this experiment, we first analyze the computational advantage of the proposed L2BRMKL and then conduct experimental comparison.

555 To facilitate the analysis, we use the commonly used SimpleMKL as a reference. In specific, we treat the computational cost of training SimpleMKL [17] on a set of samples with a preset regularisation parameter  $C$  as a unit, denoted by  $\tau_0$ . Roughly speaking, under the same setting, the computational cost of L2BRMKL and NSMKL [24] will be at the order of  $\tau_0$ , while the cost of RMKL [18] and MBMKL [19] will be higher than  $\tau_0$ .

561 To choose a suitable value for  $C$ , the four existing algorithms employ  
562 multi-fold cross-validation. Let  $k$  be the number of folds and let  $s$  be the  
563 number of candidate  $C$  values tested in the cross-validation process. Then,  
564 we can know that for these four algorithms, their cost on cross-validation  
565 will be no less than the order of  $ks\tau_0$ <sup>5</sup>. Differently, by tuning  $C$  via the  
566 MKL process, the proposed L2BRMKL can maintain the cost at the order  
567 of  $\tau_0$ , since only the number of kernel weights is slightly increased, from  $m$   
568 to  $m + 1$ . Also, among the four existing algorithms MBMKL [19] needs to  
569 compute the radius of MEB by solving a QP problem at each iteration. In  
570 contrast, the proposed L2BRMKL employs an approximation to the radius  
571 and therefore avoids such computation. This makes L2BRMKL more efficient  
572 than MBMKL in terms of integrating the radius information. Putting the  
573 above discussion together, we can see that L2BRMKL has the overall highest  
574 computational efficiency. Note that the above analysis does not depend on  
575 the machine, platform or language used to implement these algorithms.

576 As an experimental support to the above analysis, we compare the timing  
577 result of the above five MKL algorithms on three largest UCI data sets (*Ger-*  
578 *manum*, *Splice and Wdbc*) and the three AD prediction tasks. All the algo-  
579 rithms are implemented in Matlab, and no special measure is taken to opti-  
580 mize the speed of L2BRMKL. The timing result is the sum of cross-validation  
581 time and training time. The UWMKL algorithm is not included because it  
582 does not involve any learning procedure. The logarithm is applied to provide  
583 better illustration. As seen in Figure 1(a), L2BRMKL (in black) is computa-

---

<sup>5</sup>Without loss of generality, we ignore the variation on training time due to that only  $\frac{k-1}{k}$  of training samples are used in each training session of a  $k$ -fold cross-validation.

584 tionally much more efficient than the other MKL algorithms, especially when  
 585 compared with the radius-incorporated ones, RMKL (in cyan) and MBMKL  
 586 (in yellow). For example, on the data set of *Germanium*, the difference  
 587 could reach the order of  $10^3$ . Also, SimpleMKL and MSMKL generally show  
 588 similar computational cost while RMKL and MBMKL are computationally  
 589 most expensive. Figure 1(b) shows the case for three AD prediction tasks,  
 590 in which the computational advantage of the proposed L2BRMKL can still  
 be seen. These timing results are well consistent with the above analysis.

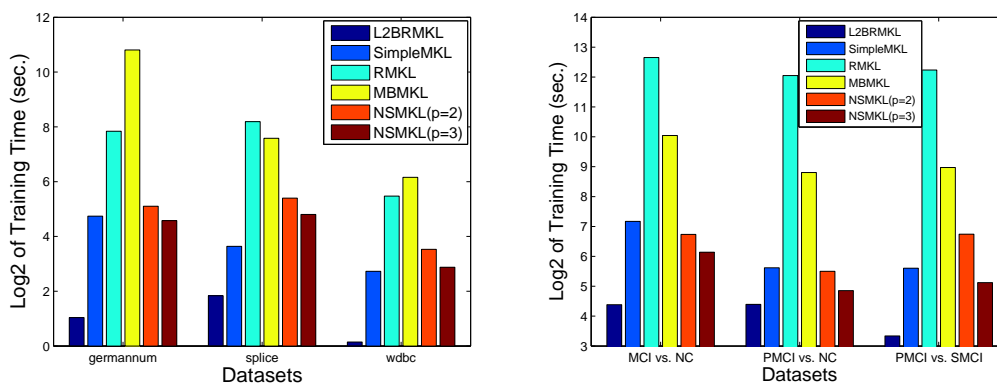


Figure 1: Experimental comparison of timing results of different MKL algorithms, where L2BRMKL is the proposed method. Figure 1(a): Timing result on three UCI data sets; Figure 1(b): Timing result on three AD prediction tasks.

591

## 592 6. Conclusion

593 Multiple kernel learning has become an effective method to predict AD  
 594 by combining information from different sources. However, the recently de-  
 595 veloped radius-incorporated MKL algorithms do not give satisfactory perfor-  
 596 mance on AD prediction tasks which are often difficult and only have a small

597 number of training samples. To improve this situation, this paper proposes  
598 an improved radius-incorporated MKL algorithm to better handle the AD  
599 prediction tasks. Instead of rigidly computing the radius of MEB, it approx-  
600 imates this radius with a linear combination of the radiuses pre-computed  
601 with each base kernel. Also, it absorbs the regularisation parameter into  
602 the MKL process and jointly optimises it with the other kernel combination  
603 weights. Through theoretical analysis, we discuss the connection of the pro-  
604 posed MKL algorithm to the well-known SimpleMKL algorithm and show  
605 that it can be readily solved. The effectiveness of the proposed algorithm  
606 is scrutinised and experimentally investigated on both pattern recognition  
607 benchmark data sets and three AD prediction tasks. As observed, it produces  
608 overall better classification performance and achieves higher computational  
609 efficiency.

610 The result in this paper also indicates that for the radius-incorporated  
611 MKL methods, how to compute the radius and design the objective function  
612 can significantly impact the kernel learning performance in practice. This  
613 raises interesting questions for both MKL research and its real applications,  
614 and they are worth exploring in the future work.

## 615 **Acknowledgement**

616 This work was supported by the National Natural Science Foundation of  
617 China (project no. 61125201, 61403405 and 60970034).

618 **Reference**

- 619 [1] P. M. Rasmussen, L. K. Hansen, K. H. Madsen, N. W. Churchill, S. C.  
620 Strother, Model sparsity and brain pattern interpretation of classifica-  
621 tion models in neuroimaging, *Pattern Recognition* 45 (6) (2012) 2085–  
622 2100.
- 623 [2] S. H. Park, S. Lee, I. D. Yun, S. U. Lee, Hierarchical MRF of glob-  
624 ally consistent localized classifiers for 3D medical image segmentation,  
625 *Pattern Recognition* 46 (9) (2013) 2408–2419.
- 626 [3] B. Caldairou, N. Passat, P. A. Habas, C. Studholme, F. Rousseau, A  
627 non-local fuzzy segmentation method: Application to brain MRI, *Pat-  
628 tern Recognition* 44 (9) (2011) 1916–1927.
- 629 [4] Alzheimers Disease and Dementia: A Comparison of International Ap-  
630 proaches, Tech. Rep., Report of the Social Committee on Aging, United  
631 States Senate, S. RES. 81, SEC. 17(d), MARCH 2, 2011.
- 632 [5] J. Bischkopf, A. Busse, M. C. Angermeyer, Mild cognitive impairment a  
633 review of prevalence, incidence and outcome according to current ap-  
634 proaches, *Acta Psychiatr Scand* 106 (2002) 403–414.
- 635 [6] J. Ye, T. Wu, J. Li, K. Chen, Machine Learning Approaches for the  
636 Neuroimaging Study of Alzheimer’s Disease, *IEEE Computer* 44 (4)  
637 (2011) 99–101.
- 638 [7] S. Duchesne, A. Caroli, C. Geroldi, C. Barillot, G. B. Frisoni, D. L.  
639 Collins, MRI-Based Automated Computer Classification of Probable AD  
640 Versus Normal Controls, *IEEE TMI* 27 (4) (2008) 509–520.



- 641 [8] J. Ye, K. Chen, T. Wu, J. Li, Z. Zhao, R. Patel, M. Bae, R. Janardan,  
642 H. Liu, G. E. Alexander, E. Reiman, Heterogeneous data fusion for  
643 alzheimer’s disease study, in: KDD, 1025–1033, 2008.
- 644 [9] H. Wang, F. Nie, H. Huang, S. Kim, K. Nho, S. L. Risacher, A. J. Saykin,  
645 L. Shen, Identifying quantitative trait loci via group-sparse multitask  
646 regression and feature selection: an imaging genetics study of the ADNI  
647 cohort, *Bioinformatics* 28 (2) (2012) 229–237.
- 648 [10] Y. Fan, S. M. Resnick, X. Wu, C. Davatzikos, Structural and functional  
649 biomarkers of prodromal Alzheimer’s disease: A high-dimensional pat-  
650 tern classification study, *NeuroImage* 41 (2) (2008) 277–285.
- 651 [11] C. Davatzikos, P. Bhatt, L. M. Shaw, K. N. Batmanghelich, J. Q. Tro-  
652 janowski, Prediction of MCI to AD conversion, via MRI, CSF biomark-  
653 ers, pattern classification, *Neurobiol Aging* .
- 654 [12] C. Hinrichs, V. Singh, G. Xu, S. C. Johnson, Predictive markers for AD  
655 in a multi-modality framework: An analysis of MCI progression in the  
656 ADNI population, *Neuroimage* 55 (2) (2011) 574–589.
- 657 [13] C. Hinrichs, V. Singh, G. Xu, S. Johnson, MKL for robust multi-  
658 modality AD classification, in: MICCAI, 786–794, 2009.
- 659 [14] D. Zhang and Y. Wang and L. Zhou and H. Yuan and D. Shen, Mul-  
660 timodal classification of Alzheimer’s disease and mild cognitive impair-  
661 ment, *NeuroImage* 55 (3) (2011) 856–867.
- 662 [15] Z. Dai, C. Yan, Z. Wang, J. Wang, M. Xia, K. Li, Y. He, Discrimi-  
663 native analysis of early Alzheimer’s disease using multi-modal imaging

- 664 and multi-level characterization with multi-classifier (M3), *NeuroImage*  
665 59 (3) (2012) 2187–2195.
- 666 [16] N. Cristianini, J. Shawe-Taylor, *An Introduction to Support Vector Ma-*  
667 *chines: And Other Kernel-based Learning Methods*, Cambridge Univer-  
668 *sity Press*, New York, NY, USA, ISBN 0-521-78019-5, 2000.
- 669 [17] A. Rakotomamonjy, F. Bach, Y. Grandvalet, S. Canu, *SimpleMKL*,  
670 *JMLR* 9 (2008) 2491–2521.
- 671 [18] H. Do, A. Kalousis, A. Woznica, M. Hilario, *Margin and Radius Based*  
672 *Multile Kernel Learning*, in: *Proceedings of the ECML*, 330–343, 2009.
- 673 [19] K. Gai, G. Chen, C. Zhang, *Learning Kernels with Radiuses of Minimum*  
674 *Enclosing Balls*, in: *NIPS*, 649–657, 2010.
- 675 [20] X. Liu, L. Wang, J. Yin, E. Zhu, J. Zhang, *An Efficient Approach to*  
676 *Integrating Radius Information into Multiple Kernel Learning*, *IEEE T.*  
677 *Cybernetics* 43 (2) (2013) 557–569.
- 678 [21] S. S. Keerthi, *Efficient tuning of SVM hyperparameters using ra-*  
679 *dius/margin bound and iterative algorithms*, *IEEE TNN* 13 (5) (2002)  
680 1225–1229.
- 681 [22] O.Chapelle, V.Vapnik, O.Bousquet, S.Mukherjee, *Choosing Multiple*  
682 *Parameters for Support Vector Machines*, *Machine Learning* 46 (2002)  
683 131–159.
- 684 [23] J. Ye and S. Ji and J. Chen, *Multi-class Discriminant Kernel Learning*  
685 *via Convex Programming*, *JMLR* 9 (2008) 719–758.

- 686 [24] Z. Xu, R. Jin, H. Yang, I. King, M. R. Lyu, Simple and Efficient Multiple  
687 Kernel Learning by Group Lasso, in: Proc. 27th ICML, 1175–1182, 2010.
- 688 [25] J. Shawe-Taylor, N. Cristianini, Kernel Methods for Pattern Analysis,  
689 Cambridge University Press, ISBN 978-0-521-81397-6, 2004.
- 690 [26] M. Gönen, E. Alpaydın, Multiple Kernel Learning Algorithms, JMLR  
691 12 (Jul) (2011) 2211–2268.
- 692 [27] S. Boyd, L. Vandenberghe, Convex Optimization, Cambridge University  
693 Press, New York, NY, USA, ISBN 0521833787, 2004.
- 694 [28] V. Vapnik, O. Chapelle, Bounds on Error Expectation for Support Vec-  
695 tor Machines, Neural Comput. 12 (9) (2000) 2013–2036, ISSN 0899-7667.
- 696 [29] L. Wang, Feature Selection with Kernel Class Separability, IEEE Trans-  
697 action on Pattern Analysis and Machine Intelligence 30 (2008) 1534–  
698 1546.