

# 1 Model selection based on combined penalties for biomarker identification

## 3 Abstract

4 The growing role of targeted medicine has led to an increased focus on the development of  
5 actionable biomarkers. Current penalized selection methods that are used to identify biomarker  
6 panels for classification in high dimensional data, however, often result in highly complex  
7 panels that need careful pruning for practical use. In the framework of regularization methods a  
8 penalty that is a weighted sum of the  $L_1$  and  $L_0$  norm has been proposed to account for the  
9 complexity of the resulting model. In practice, the limitation of this penalty is that the objective  
10 function is non-convex, non-smooth, the optimization is computationally intensive and the  
11 application to high-dimensional settings is challenging. In this paper we propose a stepwise  
12 forward variable selection method which combines the  $L_0$  with  $L_1$  or  $L_2$  norms. The penalized  
13 likelihood criterion that is used in the stepwise selection procedure results in more parsimonious  
14 models, keeping only the most relevant features. Simulation results and a real application show  
15 that our approach exhibits a comparable performance with common selection methods with  
16 respect to the prediction performance whilst minimizing the number of variables in the selected  
17 model resulting in a more parsimonious model as desired.

18  
19 **Keywords:** biomarker panels, combined penalties, model selection, penalized regression,  
20 regularization, sparsity, stepwise variable selection, treatment responder.

## 22 1. Introduction

23  
24 The high costs and long duration of clinical development, paired with high levels of attrition,  
25 require the quantification of the risk when moving from early to late stage clinical development,  
26 and biomarkers may play an important role in this quantification. However, only rarely the  
27 number of variables (biomarkers) in the resulting panel plays an active role in selection  
28 procedures. Variable selection is an important aspect in the determination of such panels in the

29 framework of high-dimensional statistical modeling. In practice, a large number of candidate  
30 predictors are available for modeling. Keeping only the relevant variables in the model makes  
31 interpretation easier and may increase the predictability of the resulting model.  
32  
33 Particularly in the framework of regularization methods, various penalty functions are used to  
34 perform variable selection. Frank and Friedman (1993) proposed the bridge regression by  
35 introducing the penalty of the form  $L_q = \sum_{j=1}^d |\beta_j|^q$ ,  $q > 0$ , for the vector of regression  
36 coefficients  $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_d) \in \mathbb{R}^d$ . When  $q \leq 1$  the penalty performs variable selection.  
37 The case where  $q = 1$  is the  $L_1$  penalty and corresponds to the Least Absolute Shrinkage and  
38 Selection Operator (Lasso) (Tibshirani, 1995). It performs continuous shrinkage and variable  
39 selection at the same time, whereas for  $q = 2$  we get the ridge estimator (Hoerl and Kennard,  
40 1970) that shrinks coefficients towards zero but it does not perform variable selection. The limit  
41 of the  $L_q$  as  $q \rightarrow 0$  gives the  $L_0$  penalty, which penalizes the number of non-zero coefficients  
42 and thus is appealing for model selection, if sparse models are of advantage. However, due to its  
43 non-convexity and discontinuity at the origin, the corresponding optimization problem becomes  
44 difficult to implement in high dimensions. In addition, the solution using  $L_0$  may be unstable  
45 because it may not be identifiable.  
46  
47 In genomic research, an  $L_1$  penalty is routinely used due to its convexity and optimization  
48 simplicity. However, the result of the  $L_1$  type regularization may not be sparse enough for a  
49 good interpretation. The development of methods to obtain sparser solutions than through  $L_1$   
50 penalization methods is becoming essential part in the classification and feature selection area.  
51 A variable selection method that combines the  $L_1$  and  $L_0$  penalties was proposed by Liu and Wu  
52 (Liu and Wu, 2007). They used a mixed integer programming algorithm for optimization of the  
53 objective function. The results showed that their method achieved sparser solutions than Lasso  
54 and more stable solutions than the  $L_0$  regularization. However the application was limited to

55 moderate data sizes, due to computational inefficiency for large-scale problems. Other  
56 combinations of  $L_q$  penalties have been proposed so far (Zou and Hastie, 2005) and recently  
57 (Huang et al., 2016) with each of these methods using a different optimization algorithm to  
58 approach the solution.

59

60 In this article, we propose a method for variable selection that penalizes the likelihood function  
61 with a linear combination of  $L_0$  with  $L_1$  or  $L_2$  penalties ( $CL$ ,  $CL2$ ) in a stepwise forward  
62 variable selection procedure. The aim is to obtain a model that is sparser than the model with  
63 the  $L_1$  penalty alone and at the same time achieve a good predictive performance. Moreover, a  
64 strong motivation for the proposed stepwise forward variable selection method is that state-of-  
65 the-art global optimization algorithms for non-smooth and nonconvex functions do not provide  
66 satisfactory results. In section 2, we define the  $CL$  and  $CL2$  penalties and present the algorithm  
67 for solving the penalized logistic regression problem with these combined penalties. In section  
68 3, we use simulated data to evaluate the performance of our method and we compare it to Lasso  
69 and adaptive Lasso both in terms of correct variable selection (true covariates with  $\beta_j \neq 0$ ) as  
70 well as predictive performance. Finally, we show an application of our method for classification  
71 and variable selection on a real dataset with protein measurements to identify the least number  
72 of predictors that can best classify responders and non-responders to a treatment.

73

## 74 **2. Methods**

75

### 76 **2.1 Regularization**

77

78 Suppose we have data  $(\mathbf{X}, \mathbf{y})$ , where  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  is the vector of responses and  $\mathbf{X}$  is an  
79  $n \times d$  matrix of predictors. We will assume that the observations are independent and the  
80 predictors standardized. With linear predictor  $\eta = \mathbf{X}^T \boldsymbol{\beta}$  and link function  $g$  the generalized  
81 linear model is expressed as

$$82 \quad g(E(\mathbf{y}|\mathbf{X})) = \eta \quad (2.1)$$

83

84 Under the regularization framework, the estimated coefficients  $\hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_d) \in \mathbb{R}^d$  are

85 obtained by minimizing the objective function  $-\log L + \lambda P(\boldsymbol{\beta})$ , and are given by:

86 
$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}}\{-\log L + \lambda P(\boldsymbol{\beta})\}$$

87 where  $P(\boldsymbol{\beta})$  is a regularization term. The parameter  $\lambda > 0$  is a tuning parameter and  $-\log L$  is the

88 negative log-likelihood. One of the most popular and commonly used regularization method is

89 the  $L_1$  regularization (Lasso), where  $P(\boldsymbol{\beta}) = \sum_{j=1}^d |\beta_j|$ . Setting  $\lambda = 0$  reverses the Lasso to

90 Maximum likelihood estimation. On the other hand, a very large  $\lambda$  will completely shrink  $\boldsymbol{\beta}$  to

91 zero thus leading to the empty or null model. In general, moderate values of  $\lambda$  will cause

92 shrinkage of the solutions towards zero, and some coefficients may be exactly zero.

93

94 Other types of  $L_1$  regularization include the adaptive Lasso, where adaptive weights are used for

95 penalizing different coefficients in the  $L_1$  penalty and was shown to have the oracle property

96 (Zou, 2006). A variable selection and estimation procedure is said to have the oracle property i)

97 if it selects the true model with probability tending to 1 and ii) if the estimated penalized

98 coefficients are asymptotically normal, with the same asymptotic empirical variance as the

99 estimator based on the true model.

100

101 However, the  $L_1$  type regularization is consistent only under rather restrictive assumptions

102 (Zhao and Yu, 2006) and the coefficient estimates are severely biased due to shrinkage

103 (Meinshausen and Yu, 2009); (Fan and Li, 2001). Although the  $L_0$  norm, where  $P(\boldsymbol{\beta}) =$

104  $\sum_{j=1}^d 1_{\beta_j \neq 0}$  and  $1_{\beta_j \neq 0}$  is the indicator function of whether  $\beta_j \neq 0$ , tend to yield the sparsest

105 solutions, its implementation in high dimensional data becomes an NP hard optimization

106 problem and is not computationally feasible. Classical information criteria like AIC (Akaike,

107 1974) or BIC (Schwarz, 1978) lie in the general class of the regularization  $\lambda P(\boldsymbol{\beta}) =$

108  $\lambda \sum_{j=1}^d 1_{\beta_j \neq 0}$  for suitable choices of  $\lambda$ . In order to gain a more concise and sparse solution and  
 109 whilst keeping a high predictive accuracy of the classification model, we propose a  
 110 regularization term that combines the  $L_0$  with  $L_1$  or  $L_2$  norms (Liu and Wu, 2007). Figure 1 plots  
 111 the penalty functions  $L_1$  and  $L_2$  in the bottom panel and the  $L_0$  penalty in the top panel. Unlike  
 112  $L_2$ , the penalty terms  $L_1$  and  $L_0$  are singular at the origin and thus perform variable selection  
 113 (Fan and Li, 2001).

114

115 **Figure 1.**

## 116 2.2 The combined $L_0 + L_1$ penalty

117

118 Following Liu and Wu (Liu and Wu, 2007) the penalization term is defined as  $CL_{\alpha}^{\varepsilon}(\beta) =$   
 119  $(1 - a)L_0^{\varepsilon} + aL_1$ , where  $0 \leq a \leq 1$  is a weighting parameter between  $L_0^{\varepsilon}$  and  $L_1$  penalties,  
 120 with  $L_0^{\varepsilon}$  given by:

$$121 \quad L_0^{\varepsilon}(\beta) = \begin{cases} 1, & |\beta| \geq \varepsilon \\ \frac{|\beta|}{\varepsilon}, & |\beta| < \varepsilon \end{cases} \quad (2.2)$$

122

123 Clearly  $CL_1^{\varepsilon} = L_1$  ( $a = 1$ ) and  $CL_0^{\varepsilon} = L_0^{\varepsilon}$  ( $a = 0$ ) are special cases of  $CL_{\alpha}^{\varepsilon}$ . Discontinuity  
 124 at the origin of  $L_0$  makes the optimization difficult and therefore we consider the continuous  
 125 approximation to  $L_0$  defined by (2.2). The limit of  $L_0^{\varepsilon}(\beta)$  when  $\varepsilon \rightarrow 0$  is  $L_0(\beta)$  itself. When  
 126  $\varepsilon > 0$  is small,  $L_0^{\varepsilon}(\beta)$  is a good approximation to  $L_0(\beta)$  (Figure 1 top right). The estimated  
 127 coefficients are obtained by minimizing the objective function

$$128 \quad -\log L + \lambda \sum_{j=1}^d CL_{\alpha}^{\varepsilon}(\beta_j) \quad (2.3)$$

129 and are given by

$$130 \quad \hat{\beta}_{CL_{\alpha}^{\varepsilon}} = \underset{\beta}{\operatorname{argmin}} \{-\log L + \lambda \sum_{j=1}^d CL_{\alpha}^{\varepsilon}(\beta_j)\}$$

131

## 132 2.3 The combined $L_0 + L_2$ penalty

133

134 We now consider another combination, the  $L_0$  norm with  $L_2$ . The motivation for combining the  
135  $L_0$  norm with  $L_2$ , is to consider a penalty that will join the nice properties of the  $L_2$  and those of  
136 the  $L_0$  norm, which is to perform variable selection ( $L_0$ ) and keep in the model groups of  
137 variables that are correlated ( $L_2$ ). In theory, a strictly convex function provides a sufficient  
138 condition for such grouping of variables and the  $L_2$  penalty guarantees strict convexity. The  
139 grouping effect refers to the simultaneous inclusion (or exclusion) of correlated predictors in the  
140 model.

141

142 The penalization term is now defined  $CL2_\alpha^\varepsilon(\beta) = (1 - a) L_0^\varepsilon + a L_2$ , where  $0 \leq a \leq 1$ . The  
143  $L_0^\varepsilon$  term introduced above is for variable selection and the  $L_2$  penalty shrinks the coefficients  
144 towards zero with no contribution to variable selection. Figure 2 gives a graphical  
145 representation of the regularization terms  $CL_{0.3}^{0.1}, L_1, L_2, CL_{0.5}^{0.1}$ .

146

147 **Figure 2.**

## 148 **2.4 The stepwise forward procedure**

149

150 In their paper Liu and Wu (2007) proposed a global algorithm to solve the corresponding  
151 difficult nonconvex problem (Mixed integer linear programming). However, the applicability  
152 was restricted to moderate datasizes. As mentioned by Frommlet F. and Nuel G. (Frommlet and  
153 Nuel, 2016), when the number of predictors grows large ( $d > 20$ ) it is not possible to apply  
154 algorithms which guarantee to find the optimal solution (Furnival and Wilson, 2000) and  
155 instead heuristic search strategies like stepwise procedures may be considered. By using  
156 heuristic techniques, we can approximate the solution of the non-smooth, non-convex and NP-  
157 hard optimization problems like the one in equation (2.3), where exact algorithms are not  
158 applicable for such minimization problems.

159

160 The optimization of the objective function (2.3) is rather challenging since  $CL_{\alpha}^{\varepsilon}(\beta)$  and  
 161  $CL2_{\alpha}^{\varepsilon}(\beta)$  are non-convex and non-differentiable at some points of the parameters' space. We  
 162 apply the BFGS, Broyden (Broyden, 1970)- Fletcher (Fletcher, 1970)- Goldfarb (Goldfarb,  
 163 1970)- Shanno (Shanno, 1970) (BFGS) variable metric (quasi Newton) method, which is shown  
 164 to work well for the optimization of non-smooth and non-convex functions (Lewis and Overton,  
 165 2009); (Lewis and Overton, 2013); (Curtis and Que, 2015). The BFGS method uses an  
 166 approximation of the Hessian matrix to find the stationary points of the function to be  
 167 minimized. Its ability to capture the curvature information of the considered function makes the  
 168 method so effective.

169  
 170 We propose to use a stepwise forward variable selection using the previously introduced  
 171 penalized likelihood criterion for feature selection that can be used effectively in high  
 172 dimensional data. In this stepwise forward selection framework, at each step we optimize the  
 173 objective function  $-\log L + \lambda a L_1$  using the BFGS algorithm and obtain

$$174 \quad \hat{\beta}_{L_1} = \underset{\beta}{\operatorname{argmin}} \{-\log L + \lambda a \sum_{j=1}^d L_1(\beta_j)\}.$$

175 The selected model is based on the criterion that minimizes the value of

$$176 \quad -\log L(\hat{\beta}_{L_1}) + \lambda a L_1(\hat{\beta}_{L_1}) + \lambda (1 - a) L_0(\hat{\beta}_{L_1}) \quad (2.4)$$

177

178 The suggested algorithm is described as follows:

179 **Step 1:**

- 180 • Given a set of  $d$  standardized predictors  $X_1, X_2, \dots, X_d$  and a response  $y_i \in \{0,1\}$ ,  $i =$   
 181  $1, \dots, n$  we consider all possible univariate models ( $M_1, M_2, \dots, M_d$ )

$$182 \quad M_1: Y \sim \beta_0 + \beta_1 X_1, \quad M_2: Y \sim \beta_0 + \beta_2 X_2, \quad M_3: Y \sim \beta_0 + \beta_3 X_3, \quad \dots, \quad M_d: Y \sim \beta_0 + \beta_d X_d$$

- 183 • Estimate  $\hat{\beta}_{L_1}^{M_1}, \dots, \hat{\beta}_{L_1}^{M_d}$  and keep  $M_j$ ,  $j \in \{1, \dots, d\}$  that gives the smallest value of the  
 184 function (2.4), e.g. keep variable  $X_2$

185

186 **Step 2:**

187 • With the model chosen in step 1 (e.g.  $M_2$ ) and in an additive way we consider all the  
188  $d - 1$  models ( $M'$ ) by adding the remaining  $d - 1$  variables one at a time to the model  $M_2$ .

189 •  $M'_1 : Y \sim \beta_0 + \beta_2 \mathbf{X}_2 + \beta_1 X_1$

190  $M'_2 : Y \sim \beta_0 + \beta_2 \mathbf{X}_2 + \beta_3 X_3$

191  $\vdots$

192  $M'_d : Y \sim \beta_0 + \beta_2 \mathbf{X}_2 + \beta_d X_d$

193

194 • Keep the model that minimizes the function in (2.4),

195

196 **Step 3:**

197 • Continue adding single variables until the value of the function (2.4) in the current step  
198 is bigger than its value in the previous step.

199

200 The advantage of using the function (2.4) instead of (2.3) in the optimization is that we no

201 longer need to consider the continuous approximations to the discontinuous  $L_0$  function and

202 therefore we eliminate the number of parameters by the continuity parameter  $\varepsilon$ . The reason why

203 we can do so is that within each step the  $L_0$ -penalty term remains constant (since the dimension

204 of the model is fixed) and hence play no role in the determination of the regression coefficients.

205 The  $L_0$ -penalty term does only play a role for the stopping criterium.

206

207 **2.5 Sparse logistic regression with combined penalties**

208

209 As a particular example we consider the binary linear regression model (2.1), where  $y \in \{0,1\}$ ,

210 is a vector of  $n$  observed binary outcomes,  $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_d) \in \mathbb{R}^d$  is the vector of

211 coefficients. The link function is the logit function  $logit(p) = \log\left(\frac{p}{1-p}\right)$ , where  $p$  is the

212 conditional event probability and is given by



213 
$$p = P(\mathbf{y} = 1|\mathbf{X}) = \frac{e^\eta}{1+e^\eta}$$

214 (2.5)

215 The coefficient estimates are obtained by minimizing (2.3) with the log-Likelihood

216 
$$\log L = L(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) = \sum_{i=1}^n y_i \log(p_i) + (1 - y_i)\log(1 - p_i)$$

217

218 **3. Results**

219

220 In this section we examine via simulations the performance of logistic regression when models

221 are selected and estimated with the above introduced combined penalties ( $CL$ ,  $CL2$ ) by either

222 stepwise forward selection as introduced in Section 2 (*stepCL* and *stepCL2*) or by global

223 minimization ( $CL$ ,  $CL2$ ). In the stepwise model selection scheme, we also examine the

224 performance of the stepwise adaptive  $L_1$  model with the  $\lambda (1 - a) L_0$  selection criterion

225 (*stepAdaCL*). For that model the objective function to minimize is  $-\log L + \lambda a \sum_{j=1}^d w_j L_1(\beta_j)$ ,

226 where  $w_j = \frac{1}{|\beta_j^*|^{\nu}}$  are the adaptive weights and  $|\beta_j^*|$  is the ridge regression estimator. The

227 estimation is done with the stepwise algorithm described in Section 2.4.

228

229 **Although the proposed method is a stepwise variable selection procedure, we did not consider**

230 **the comparison with other stepwise methods like the stepwise BIC or AIC, as they tend to**

231 **perform poorly when the dimension is large relative to the sample size and are usually too**

232 **liberal, that is, they tend to select a model with many spurious covariates (Chen and Chen,**

233 **2008). As mentioned by Zhang and Shen (Zhang and Shen, 2010) these criteria may be**

234 **inadequate due to their nonadaptivity to the model space and infeasibility of exhaustive search.**

235

236 We include the global minimization in spite of the disadvantages mentioned in Section 2 for a

237 comparison. We also consider the results from Lasso ( $L_1$  penalty) and the adaptive Lasso. We

238 compare the different methods in terms of the fraction of correctly selected variables and the

239 prediction classification accuracy. The real data come from a biomarker study concerned with  
240 protein measurements with the objective to select biomarkers that potentially discriminate  
241 between responders and non-responders.

242

### 243 **3.1 Simulation Study**

244

245 We simulate data for varying number of predictors. We consider two settings, one high  
246 dimensional data where the number of predictors ( $d$ ) exceed the number of samples ( $n$ ), and a  
247 setting where the sample size is smaller than the dimensionality of the data. We assume  
248 multivariate normal predictors  $X_1, \dots, X_d$  with pairwise correlation  $\rho$  (compound symmetry). Let  
249  $\rho$  denote the correlation between variables  $X_m, X_l$  where  $m, l \in \{1, \dots, d\}, m \neq l$ .

250

251 The true model that was used to generate the outcome has  $k$  informative covariates  $X_k, k \in \mathbb{Z},$   
252  $1 < k < d$ . We consider a classification problem with  $y$  a binary response and standard  
253 normally distributed predictors  $\mathbf{X} \sim MVN(0, \Sigma)$ , where  $\Sigma$  is the covariance matrix. We consider  
254 the logistic model with logit link function,  $logit(p) = X^T \beta$ , as described above with  $p$  the  
255 probability of  $y=1$  given  $X$  as defined in (2.5). In other words, each component of the response  
256 vector  $y$  is viewed as a realization of a Bernoulli random variable with probability of success  $p$ .

257

258 Four scenarios will be presented here, each with  $n=100$  samples.

259

#### 260 1. Scenario 1: $d < n$ , correlation $\rho = 0.5$

261 We consider  $d=50$  covariates, with  $k=3$  informative predictors and coefficient vector  $\beta =$   
262  $(3, 1.5, 2, \underbrace{0, \dots, 0}_{47})$ . The correlation between the 3 informative predictors is  $\rho = corr(X_m, X_l) =$   
263  $0.5, m \neq l$  and  $m, l = 1, 2, 3$ .

264

#### 265 2. Scenario 2: High dimensional setting $d > n$ , correlation $\rho = 0.5$

266 We consider  $d=150$  covariates, with  $k=15$  informative predictors with  $\beta =$   
 267  $(\underbrace{3, \dots, -3.5, \dots}_3, \underbrace{1.5, \dots, 5, \dots}_3, \underbrace{-2, \dots, 0, \dots, 0}_{135})$ . The correlation  $\rho$  between the 15 informative  
 268 predictors is  $\text{corr}(X_m, X_l) = 0.5$ ,  $m \neq l$  and  $m, l = 1, \dots, 15$ .

269

270 3. Scenario 3: High dimensional setting  $d > n$ , correlation  $\rho = 0.7$

271 The dataset consists of  $n=100$  samples and  $d=200$  covariates, with  $k=15$  informative predictors  
 272 with  $\beta = (\underbrace{2, \dots, -3, \dots}_4, \underbrace{1.5, \dots, -2, \dots}_4, \underbrace{0, \dots, 0}_{185})$ . The correlation  $\rho$  between the 15 informative  
 273 predictors is  $\text{corr}(X_m, X_l) = 0.7$ ,  $m \neq l$  and  $m, l = 1, \dots, 15$

274

275 4. Scenario 4: High dimensional setting  $d > n$ , block correlation

276 We consider  $d=200$  covariates, with  $k=16$  informative predictors with  $\beta =$   
 277  $(1, \underbrace{4, \dots, -3, \dots}_4, \underbrace{1.5, \dots, -2, \dots}_4, \underbrace{0, \dots, 0}_{184})$ . In this scenario there are two groups (blocks) of  
 278 correlated predictors and one single independent feature. The coefficients of  $d-k=184$  variables  
 279 were set to zero,  $\beta_r = 0$ ,  $r = 184, \dots, 200$ . The correlation between predictors in block 1 is  
 280  $\text{corr}(X_m, X_l) = 0.4$ ,  $m \neq l$  and  $m, l = 1, \dots, 7$  and the correlation among predictors in block 2 is  
 281  $\text{corr}(X_m, X_l) = 0.7$ ,  $m \neq l$  and  $m, l = 8, \dots, 16$ .

282

283 **3.2 Tuning of parameters**

284

285 All analyses were done in R version 3.2.3 ([R Core Team, 2015](#)). For the Lasso and Adaptive  
 286 Lasso the *glmnet* library was used and all the functions that were used for the combined penalty  
 287 approach can be found in the R-package “*stepPenal*”, available on the CRAN. For the adaptive  
 288 lasso weights were estimated by ridge regression and then used for a weighted  $L_1$  penalty in  
 289 estimation of  $\beta$ . The optimal regularization parameters for the methods *stepCL*, *stepAdaCL*, *CL*,  
 290 *CL2*, *stepCL2* were tuned by 10-fold cross-validation on the two dimensional surface  $(a, \lambda)$   
 291 using a grid of values. The choice of the optimal parameters was done in the following way. For

292 each configuration of  $(a, \lambda)$  in the grid, the AUC of the ROC curves on the validation set was  
293 computed in each of the 10 validation sets. The average of the 10 AUCs was reported together  
294 with its standard deviation.

295

296 Selection of  $(a, \lambda)$  was based on the interval  $A = [\max AUC - sdAUC, \max AUC)$  where  
297  $\max AUC$  is the maximum average AUC and  $sdAUC$  is the standard deviation of the AUCs  
298 corresponding to the  $(a, \lambda)$  with maximum average AUC. The  $(a, \lambda)$  that corresponds to the  
299 median of the AUCs in the interval  $A$  was chosen for the final model fitting. In case that more  
300 than one configurations yields the median of the AUCs, we select the configuration with the  
301 largest  $\lambda$  and smallest  $a$ , to obtain the sparsest model. The use of the interval  $A$  acknowledges  
302 the sample variability and the fact that we are aiming for a compromise between good  
303 classification performance and complexity of the model. In other words, a small decrease in the  
304 AUC of the ROC curve is acceptable in return to a less complex model. The Lasso and adaptive  
305 lasso were also tuned by 10-fold cross-validation on the one dimensional space ( $\lambda$ ), using the  
306 default settings in R in the function `cv.glmnet` and the measure type "auc".

307

### 308 **3.3 Simulation Results**

309

310 The different classifiers were built by the estimated tuning parameters on the training set. Then,  
311 the obtained classifiers were applied to the testing set for classification and prediction. For the  
312 testing set, we simulated data from the same distribution as the training set for  $n=1000$  samples.  
313 **We simulated 1000 datasets on which we applied all methods.** For each method we computed  
314 the mean classification performance of the models on the testing sets measured by the AUC of  
315 the ROC curve (test AUC). This is a measure for the discrimination ability of the model to  
316 correctly distinguish the two classes of the response. The complexity of the resulting model was  
317 measured by the ratio of correctly selected variables (true covariates with  $\beta_j \neq 0$ ) to the total  
318 variables selected by the model. We will call this ratio  $RCV$  for the rest of the paper.

319

320 This ratio takes values between zero and one. When the model selects none of the informative  
321 variables it is zero and it becomes one, when the selected model includes only the  $k$  informative  
322 covariates. The closer the  $RCV$  is to one, the sparser the model is and selects the true variables.  
323 The results in Table 1 summarize the performance of the different methods in terms of model  
324 complexity. An ideal model selection method would only select the  $k$  true features and set the  
325 coefficients of the other  $d-k$  variables equal to zero.

326

327 **Table 1:**

328

329 In most of the scenarios, the *stepCL* and *stepCL2* methods yield a higher  $RCV$  than the other  
330 methods and on average the *stepCL* and *stepCL2* models are sparser than the other methods. In  
331 scenario 2 and 3 the adaptive Lasso yields the higher  $RCV$ , but the models are not as sparse as  
332 the stepwise methods. Although the stepwise methods (*stepCL*, *stepCL2*, *stepAdaCL*) result in  
333 including the least variables in the model, its discriminative ability in terms of AUC, as shown  
334 in Table 2, is comparable with the other methods that tend to select a larger model with more  
335 variables.

336

337 The *stepCL2* method also has remarkable performance both in terms of sparsity and predictive  
338 discrimination. Considering the trade-off between model complexity and performance, the  
339 proposed stepwise combined penalty approach achieves a good balance between parsimony,  
340 including less variables and maintaining a high predictive accuracy that is comparable with  
341 state-of-the-art methods. We should mention that in scenarios 2,3 and 4 none of the methods  
342 select all of the  $k$  informative variables, however, for the stepwise method the AUC of the ROC  
343 curve on the testing set is greater than 90%, indicating a good discrimination accuracy by  
344 including the least variables in the model. In all scenarios, we found that the *stepCL2* method  
345 has comparable performance to *stepCL* and is superior to adaptive Lasso and Lasso.

346

347 In Table 2 we present results regarding the predictive classification accuracy of the methods by  
348 the AUC of the ROC curves. We report the Brier score (Brier, 1950) as a measure of the  
349 accuracy of predictions, defined as

$$350 \quad Brier = \frac{1}{n} \sum_{i=1}^n (\hat{p}_i - y_i)^2$$

351 It is given by the squared distance between the patients observed status  $y_i$  and the predicted  
352 probability  $\hat{p}_i$ . The decision space for the Brier score is the interval  $[0,1]$  and generally the  
353 lowest the Brier score, the better the classification rule. If the predicted probability is 0.5 for  
354 each individual, the Brier score of 0.25 would indicate that the classification rule is a random  
355 one.

356

357 **Table 2:**

358

359 Empirical results from our simulations show that even for no high-dimensional settings where  
360  $n > d$ , the stepwise method gives the sparsest solutions whilst maintaining classification  
361 performance measures as good as Lasso and adaptive Lasso. The *CL2* method tends to select  
362 big models, due to the  $L_2$  norm which shrinks coefficients towards zero without variable  
363 selection. Thus when  $a$  is close to 1, the model will behave similar to ridge regression and the  
364 resulting model will be complex in terms of the number of predictors. On the other hand, when  
365  $a$  is closer to 0, the *CL2* and *stepCL2* penalties will borrow more of the characteristics of the  $L_0$   
366 norm and will result in sparser models.

367

368 In our simulations we also considered the case where there is no correlation among predictors  
369 (results not shown in the table as we don't consider it a realistic scenario). We repeated scenario  
370 2 with the only alteration of setting  $\rho = 0$ . Results were in the same direction as in scenario 2  
371 shown in Table 1 and Table 2. That is, the stepwise methods perform better than all the other

372 methods in terms of model complexity resulting in the sparsest models with a high classification  
373 performance.

374

375 Furthermore, we examined the situation where there are no predictors in the data associated  
376 with the outcome. In that case that the true model is the null model, none of the methods  
377 identified the true model. Again, running through again the second scenario with  $d=150$   
378 predictors with none being informative for the outcome, the *stepCL* method selected a median  
379 of 5 variables whereas the other methods selected between 14 (*AdaLasso*, *CL*) and 19 (*CL2*).

380 We observed the same pattern in the results for repeating the first scenario with  $d=50$   
381 uninformative predictors, where none of the methods selected the true model but the stepwise  
382 methods produced the sparsest solutions.

383

#### 384 **3.4 Application- real data analysis**

385

386 To illustrate the applicability of the proposed method, we applied the stepwise method on a real  
387 example involving protein measurements. The dataset contained  $n=53$  patients with baseline  
388 measurements of  $d=187$  proteins. To maintain confidentiality, the names of the proteins are not  
389 revealed. For the presentation of the results and keeping the study anonymized we renamed the  
390 proteins to  $X_1, X_2, \dots, X_{187}$ . The objective is to extract potential candidate markers  
391 discriminating responders from non-responders based on patients' protein levels. We apply our  
392 proposed stepwise combined penalty approach with the aim to select a small set of proteins that  
393 can sufficiently predict response to the treatment. We compare our approach with the commonly  
394 used Lasso and adaptive Lasso, but also with the global optimization penalized methods *CL* and  
395 *CL2*.

396

397 The regularization parameters were not tuned by cross-validation, due to the relatively small  
398 sample size ( $n=53$ ). The tuning was done using the bootstrap method in the following way; for a

399 grid of values of  $(\alpha, \lambda)$ , we trained the models on  $B=100$  bootstrapped datasets (drawing  
400 samples with replacement from the original data) and evaluate their classification performance  
401 (in terms of AUC) on the original data. For each combination of the tuning parameters  $(\alpha, \lambda)$ ,  
402 the models were trained on  $B$  bootstrapped sets and validated on the testing set (original data)  
403 and the average AUC (over the  $B$  bootstrapped samples) was reported together with its standard  
404 deviation. The configuration of  $(\alpha, \lambda)$  that corresponds to the median of the AUC in interval  $A$ ,  
405 as described above in the section 3.2 ‘Tuning of parameters’, was chosen.

406

407 The results show that the stepwise methods yield the sparsest models by selecting 8 variables  
408 (*stepCL*) and 9 (*stepCL2*) accordingly, whereas the other methods select between 22 (*CL2*) and  
409 26 (*Lasso*). It is noticeable that the classification performance of the stepwise method is as good  
410 as the other variable selection methods, albeit including the least predictors. In order to evaluate  
411 the performance of the models and in the absence of an external validation dataset we use  
412 bootstrapping. We applied all the methods on another  $B=1000$  bootstrapped datasets of the  
413 protein data, by sampling with replacement, and the frequencies of the top 10 selected variables  
414 by all methods are reported in Figure 3. This results in 16 unique proteins.

415

416 This figure shows that the proteins that were frequently selected by the *stepCL* and *stepCL2*  
417 methods are also frequently selected by the Lasso and adaptive Lasso. Note that the stepwise  
418 methods have lower frequencies of the selected variables, because selection of larger models  
419 will automatically increase the number of selection for individual variables.

420

421 **Figure 3:**

422

423 Figure 4 shows boxplots of the total number of variables included in the model over the  
424 bootstrap evaluations. The stepwise method yields consistent model selection by selecting a



425 median of 8 variables for *stepCL* and *stepCL2*, whereas the Lasso and Adaptive Lasso have a  
426 big variability on the complexity of the model selected. The AUC of the ROC curves that is  
427 used as a measure of classification performance of the methods on the bootstrapped datasets and  
428 their distribution is shown in Figure 5. The stepwise methods tend to always select the most  
429 sparse models more systematically, while maintaining a very good classification performance.

430

431 **Figure 4:**

432

433 **Figure 5:**

434

#### 435 **4. Conclusion**

436

437 In this paper we have proposed a stepwise forward approach for model selection in the  
438 framework of penalized regression using a penalty that combines the  $L_0$  norm, which is based  
439 on the number of coefficients, with  $L_1$  norm which is based on the size of coefficients or  
440  $L_2$  norm which take into account the grouping effect. The aim of the proposed method is to find  
441 a model that includes as less and relevant variables on one hand, and have good predictive  
442 performance on the other hand. The combined penalization term  $CL_\alpha^\varepsilon(\beta)$  that was introduced  
443 by Liu and Wu (2007) was limited to moderate datasets due to limitations of the optimization  
444 algorithm. Considering the heuristic stepwise forward approach, we can apply the penalization  
445  $CL_\alpha(\beta)$  and  $CL2_\alpha(\beta)$  to high-dimensional data by using the BFGS algorithm which is found  
446 to work well in practice for nonconvex and non-smooth functions (Lewis and Overton, 2009);  
447 (Lewis and Overton, 2013); (Curtis and Que, 2015). As a result, the practical implementation of  
448 the stepwise penalization method is simpler and more efficient.

449

450 We found that for the stepwise method the computational time was shorter than the global  
451 optimization. However, the tuning of the regularization parameters  $(a, \lambda)$  can be  
452 computationally intensive. **This is an important aspect of penalization methods and can be**

453 further explored and improved in future work. Simulation results and a real data application  
454 show that the proposed method yields sparser models, while maintaining a good classification  
455 performance. This is an important consideration for classification and screening applications  
456 where the goal is to develop a test using as less features as possible to control the cost. Overall,  
457 we found that our method provides a sparser model whilst maintaining similar prediction  
458 properties with the other methods. We hope that this paper could be a first step to learn more  
459 about the theoretical properties of this method, which seems to be worth of further investigation.

460

461 Furthermore, it would be of great interest to extend the forward stepwise method to the stepwise  
462 bidirectional approach, considering at each step of the algorithm which variables can be  
463 included and excluded (forward and backwards variable selection) in the model. As future work  
464 we also consider to apply our method to regression problems for variable selection with a  
465 continuous response as well as time-to-event data.

466

#### 467 **Declaration of Conflicting Interests**

468

469 The author(s) declared no potential conflicts of interest with respect to the research, authorship,  
470 and/or publication of this article.

471

#### 472 **Funding**

473

474 This project has received funding from the European Union's Horizon 2020 research and  
475 innovation programme under the Marie Skłodowska-Curie grant agreement No 633567 and is  
476 part of the IDEAS European training network (<http://www.ideas-itn.eu/>). This report is in part  
477 independent research arising from Prof Jaki's Senior Research Fellowship (NIHR-SRF-2015-  
478 08-001) supported by the National Institute for Health Research. The views expressed in this  
479 publication are those of the authors and not necessarily those of the NHS, the National Institute  
480 for Health Research or the Department of Health.

481 **References**

482

483 Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on*  
484 *Automatic Control*, p. 19(6): 716–723.

485 Brier, G., 1950. Verification of forecasts expressed in terms of probability. *Mon Wea Rev.*, pp.  
486 78;1-3.

487 Broyden, C. G., 1970. The convergence of a class of double-rank minimization algorithms..  
488 *IMA Journal of Applied Mathematics* 6(1), pp. 76-90.

489 Chen, J. and Chen, Z., 2008. Extended Bayesian information criteria for model selection with  
490 large model spaces. *Biometrika*, pp.759-771.

491 Curtis, F.E. and Que, X., 2015. A quasi-Newton algorithm for nonconvex, nonsmooth  
492 optimization with global convergence guarantees. *Mathematical Programming Computation*,  
493 7(4), pp.399-428.

494

495 Fan, J. and Li, R., 2001. Variable selection via nonconcave penalized likelihood and its oracle  
496 properties. *Journal of the American statistical Association*, 96(456), pp.1348-1360.

497

498 Fletcher, R., 1970. A new approach to variable metric algorithms. *Compt.J* 13(3), pp. 317-322.

499 Frank, L.E. and Friedman, J.H., 1993. A statistical view of some chemometrics regression tools.  
500 *Technometrics*, 35(2), pp.109-135.

501

502 Frommlet, Florian, and Gregory Nuel. "An Adaptive Ridge Procedure for L 0 Regularization."  
503 *PloS one* 11, no. 2 (2016): e0148620.

504

505 Furnival, George M., and Robert W. Wilson. "Regressions by leaps and bounds." *Technometrics*  
506 16, no. 4 (1974): 499-511.

507

508 Goldfarb, D., 1970. A family of variable-metric methods derived by variational means.  
509 *Mathematics of computation* 24(109), pp. 23-26.

510 Hoerl, A.E. and Kennard, R.W., 1970. Ridge regression: Biased estimation for nonorthogonal  
511 problems. *Technometrics*, 12(1), pp.55-67.

512

513 Huang, H.H., Liu, X.Y. and Liang, Y., 2016. Feature Selection and Cancer Classification via  
514 Sparse Logistic Regression with the Hybrid L 1/2+ 2 Regularization. *PloS one*, 11(5),  
515 p.e0149675.

516

517 Lewis, Adrian S., and Michael L. Overton. "Nonsmooth optimization via BFGS." *Submitted to*  
518 *SIAM J. Optimiz* (2009): 1-35.

519

520 Lewis, A.S. and Overton, M.L., 2013. Nonsmooth optimization via quasi-Newton methods.  
521 *Mathematical Programming*, pp.1-29.

522  
523 Liu, Y. and Wu, Y., 2007. Variable selection via a combination of the L 0 and L 1 penalties.  
524 *Journal of Computational and Graphical Statistics*, 16(4), pp.782-798.  
525  
526 Meinshausen, N. and Yu, B., 2009. Lasso-type recovery of sparse representations for high-  
527 dimensional data. *The Annals of Statistics*, pp.246-270.  
528  
529 R Core Team, 2015. *R: A Language and Environment for Statistical Computing*, Vienna,  
530 Austria: R Foundation for Statistical Computing.

531 Schwarz, G. E., 1978. Estimating the dimension of a model. *Annals of Statistics*, p. 6 (2): 461–  
532 464.

533 Shanno, D. F., 1970. Conditioning of quasi-Newton methods for function minimization.  
534 *Mathematics of computation*, 24(111), pp. 647-656.

535 Tibshirani, R., 1996. Regression Shrinkage and Selection via the Lasso. *J.R. Statist. Soc. B*, pp.  
536 267-288.

537 Zhang, Y. and Shen, X., 2010. Model selection procedure for high-dimensional data. *Statistical*  
538 *analysis and data mining*, 3(5), pp.350-358

539 Zhao, P. and Yu, B., 2006. On model selection consistency of Lasso. *Journal of Machine*  
540 *learning research*, 7(Nov), pp.2541-2563.  
541  
542 Zou, H., 2006. The adaptive lasso and its oracle properties. *J Am Stat Assoc. Taylor & Francis*,  
543 pp. 101:1418-1429.

544 Zou, H. and Hastie, T., 2005. Regularization and variable selection via the elastic net. *Journal*  
545 *of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2), pp.301-320.  
546