

Selection Stability in High Dimensional Statistical Modelling: Defining a Threshold for Robust Model Inference

Martin Green (✉ martin.green@nottingham.ac.uk)

University of Nottingham

Eliana Lima

European Food Safety Authority

Robert Hyde

University of Nottingham

Research Article

Keywords: Covariate selection, selection stability, stability threshold, high dimensional data, statistical triangulation.

Posted Date: August 2nd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-738092/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Title**

2 Selection stability in high dimensional statistical modelling: Defining a threshold for robust
3 model inference

4

5 **Author names and affiliations**

6 Martin Green^{a*}, Eliana Lima^{a,b}, Robert Hyde^a

7

8 ^a School of Veterinary Medicine and Science, University of Nottingham, Sutton Bonington Campus,
9 Leicestershire, United Kingdom

10 ^b Present address: European Food Safety Authority, Via Carlo Magno 1A, 43126, Parma, Italy

11

12 ***Corresponding author**

13 Martin Green

14 E-mail: martin.green@nottingham.ac.uk

15 School of Veterinary Medicine and Science, University of Nottingham, Sutton Bonington Campus,
16 Leicestershire, LE12 5RD, United Kingdom

17 Tel +44 115 951 6116

18

19

20 Declarations of interest: none

21

22 **Abstract**

23 Epidemiological research commonly involves identification of causal factors from within
24 high dimensional (wide) data, where predictor variables outnumber observations. In this
25 situation, however, conventional stepwise selection procedures perform poorly. Selection
26 stability is one method to aid robust variable selection, by refitting a model to repeated
27 resamples of the data and calculating the proportion of times each covariate is selected. A key
28 problem when applying selection stability is to determine a threshold of stability above which
29 a covariate is deemed ‘important’.

30 In this research we describe and illustrate a two-step process to implement a stability
31 threshold for covariate selection. Firstly, covariate stability distributions were established
32 with a permuted model (randomly reordering the outcome to sever the relationship with
33 predictors) using a cumulative distribution function. Subsequently, covariate stability was
34 estimated using the true model outcome and covariates with a stability above a threshold
35 defined from the permuted model, were selected in a final model. The proposed method
36 performed well across 22 varied, simulated datasets with known outcomes; selection error
37 rates were consistently lower than conventional implementation of equivalent models. This
38 method of covariate selection appears to offer substantial advantages over current methods, to
39 accurately identify the correct covariates from within a large, complex parameter space.

40

41

42

43 **Keywords**

44 Covariate selection; selection stability; stability threshold; high dimensional data; statistical
45 triangulation.

46

47 **1. Introduction**

48 Epidemiological research increasingly involves identification of potentially causal factors
49 from within relatively wide datasets, that is, when the number of variables (p) is relatively
50 large compared to the number of observations (n). Identification of causal variables for
51 inference from within a high dimensional, wide dataspace is problematic because
52 conventional stepwise selection procedures perform poorly, resulting in over fit models ¹⁻⁴.
53 Statistical methods have been developed to improve variable selection with such data,
54 including modifications to AIC/BIC ⁵, and a range of regularisation methods based on
55 functions that penalise model coefficients to balance over and under fitting, the so called
56 variance-bias trade off ⁶⁻⁸.

57 In addition to regularisation, it is acknowledged that robustness in model selection is
58 improved through use of covariate selection stability ⁹⁻¹¹. Covariate stability is estimated by
59 refitting a model to repeated resamples of the data (for example through Bootstrapping) and
60 calculating the proportion of times each covariate is selected across all resamples. Whilst it is
61 known that the most stable variables (those selected in most subsampled models) are least
62 likely to be false positives ⁹, the optimal threshold for stability, above which a covariate
63 should be deemed 'important' or 'significant', has not been determined. A stability threshold
64 originally proposed for use with lasso regression ⁹ has been shown to be too conservative,
65 resulting in true causal variables being missed ^{12,13}. An empirical method of stability
66 selection, proposed for genetic data using elastic net regression ¹³, appeared to improve upon
67 the method proposed by Meinshausen and Bühlmann ⁹, but the issue of missing many true
68 causal variables (false negative results) remained. To date, since a clear, generalisable
69 method to identify a stability threshold is unavailable, arbitrary thresholds have been
70 employed for stability analyses in veterinary epidemiology ^{14,15} and the need for research to

71 establish a suitable cut-off for the covariate selection in high dimensional data is clear and
72 has been recently re-emphasised in human epidemiology ¹⁶.
73 In this paper we build on principles described by Meinshausen and Bühlmann ⁹ and Kim et al
74 ¹³, to develop a general approach to identify a threshold for covariate selection stability and
75 hence to aid robust model selection for high dimensional, correlated data. To evaluate and
76 illustrate the method, we use 22 simulated datasets with known outcomes and increasingly
77 complex correlation structures, alongside five methods of covariate selection to illustrate
78 generalisability. For each statistical method we compare the performance of covariate
79 selection based on a conventional approach to that using the application of selection stability
80 with the newly defined threshold.

81

82 2. Materials and Methods

83 Five statistical methods were used to analyse 22 simulated datasets that were constructed
84 with known underlying relationships and correlation structures. For each method, a
85 conventional approach to covariate selection was conducted followed by implementation of
86 selection stability by bootstrapping ¹⁴. The modelling methods chosen are acknowledged
87 approaches for inferential modelling and were stepwise selection based on Akaike
88 information criterion (sAIC) ¹⁷, stepwise selection based on a modified Bayesian Information
89 Criterion (mBIC) ⁵, elastic net regression (enet) ⁸, minimax convex penalty regression (MCP)
90 ¹⁸ and a combination method that synthesised results from all four methods ¹⁹. These
91 modelling techniques are described in detail in Section 2.1 and the new approach to stability
92 selection, including the determination of a threshold for inference, is described in Section 2.2.
93 All datasets used for analysis are described in Section 2.3.

94

95 2.1. Statistical models

96 The following statistical approaches were used to evaluate and compare covariate selection.

97 2.1.1. Stepwise selection using AIC (sAIC)

98 A conventional linear regression model was implemented using the ‘stepwise’ function in the
99 bigstep package²⁰ in R²¹. The regression equation took the form;

100
$$y = \beta_0 + \sum_{j=1}^p \beta_j x_j + e \quad (1)$$

101 where y was the response variable specified in the simulated data, β_0 an intercept term, x_j
102 represented the j^{th} of p covariates with an estimated coefficient β_j and e the residual model
103 error. To speed computation, variable selection was conducted by first removing explanatory
104 variables with a relatively poor correlation with the outcome (Pearson correlation, $P > 0.05$)
105 followed by an iterative stepwise procedure (forward and backward) with selection of the
106 final set of variables through minimisation of the Akaike information criterion (AIC). The
107 AIC loss function is defined as; $2k - 2 \ln(\hat{L})$, where k is the number of parameters in the
108 model and \hat{L} the likelihood function.

109

110 2.1.2. Stepwise selection using a modified Bayesian Information Criterion (mBIC)

111 A modified Bayesian Information Criterion was used for model selection, also implemented
112 in the bigstep package²⁰ with the ‘stepwise’ function being used for covariate selection. The
113 modified BIC increases the penalty on the number of parameters selected beyond the
114 conventional BIC, producing a sparser model. The mBIC loss function was implemented as;

115

116
$$mBIC = \log L(Y|M_i, \theta_i) - \frac{1}{2} k_i \log n - k_i \log\left(\frac{1-p}{p}\right) \quad (2)$$

117

118 where $\log L(Y|M_i, \theta_i)$ represented the log likelihood given model M_i and parameter

119 values θ_i , k_i was the number of predictors in the selected model, n the sample size, and p

120 the probability that a randomly chosen predictor influenced Y. As the number of available
 121 predictors increases relative to the number of samples (n), p decreases and $k_i \log(\frac{1-p}{p})$
 122 becomes of increasing importance as the penalty term ⁵.

123

124 2.1.3. Elastic net regression (enet)

125 Elastic net is a form of regularised regression that incorporates a mixture of lasso (L1) and
 126 ridge (L2) penalties ⁸ and was implemented as;

127

$$128 \quad SSE_{enet} = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda_E [\sum_{j=1}^p \frac{1}{2} (1 - \alpha) \beta_j^2 + \alpha |\beta_j|] \quad (3)$$

129

130 where SSE_{enet} represented the elastic net loss function to be minimised, i denoted each
 131 observation and n the sample size, y_i and \hat{y}_i were the simulated true and model predicted
 132 outcomes respectively for the i th observation, j denoted a predictor variable with p the
 133 number of predictor variables in total, and $|\beta|$ represented absolute values of the
 134 regression coefficients. The hyperparameters that provide the penalty (λ_E) and the relative
 135 proportion of penalisation on either the sum of the square of the coefficients or the
 136 unsquared coefficients (α) were optimised using 10 x 10-fold cross validation to minimise
 137 mean absolute error (MAE). Elastic net models were built using the glmnet package ²²
 138 using the caret package platform ²³ in R ²¹.

139 2.1.4. Minimax convex penalty (MCP)

140 Minimax convex penalty (MCP) ¹⁸ is a form of regularised regression in which the size of the
 141 penalty function varies with the size of variable coefficient, β . It was implemented as;

142

$$143 \quad SSE_{mcp} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{j=1}^p P(\beta_j | \lambda, \gamma) \quad (4)$$

144

145 where SSE_{mcp} represented the loss function to be minimised, i , y_i , \hat{y}_i , j , p and n were as
146 defined in equation (3) and $P(\beta_j|\lambda, \gamma)$ represented a penalty function as follows;

147

$$P(\beta|\lambda, \gamma) = \begin{cases} \lambda|\beta| - \frac{\beta^2}{2\gamma}, & \text{if } |\beta| \leq \gamma\lambda \\ 0.5\gamma\lambda^2, & \text{if } |\beta| > \gamma\lambda \end{cases} \quad (5)$$

150

151 where γ and λ were hyperparameters optimised using 10 x 10-fold cross validation to
152 minimise the MAE. MCP models were estimated using the R package `ncvreg` ²⁴.

153

154 2.2. Estimation of covariate selection stability and coefficient distributions

155 Conventional bootstrapping ²⁵ was used to estimate covariate stability for all analytic
156 approaches, according to methods previously described ²⁶. In brief, selection stability ^{9,10,27}
157 was evaluated for each model as the percentage of times each covariate was selected across
158 200 bootstrap samples. For elastic net and MCP, model hyperparameters used in the
159 bootstrap sampling were those identified using 10 x10-fold cross validation in an initial
160 model using the full dataset. The distributions of variable coefficients were calculated from
161 all non-zero values of the coefficient in the bootstrap samples and a bootstrap P value (BPV)
162 was calculated for each covariate ¹⁴. The BPV was calculated from all non-zero coefficient
163 values of the bootstrap sample (i.e. when the variable was selected in the model) and defined
164 as the smallest proportion of (non-zero) coefficient values to one side of zero. For example, if
165 a covariate was selected in the model in 100 bootstrap samples and 95 of these were either
166 greater or less than zero, then the Bootstrap P value would be; $(100-95)/100 = 0.05$. Variable
167 selection was visualised by plotting covariate selection stability against Bootstrap P value.

168

169 2.3. Multiple method synthesis

170 Covariate selection stabilities were combined across all modelling methods using an
171 approach described previously¹⁹. In brief, the bootstrap matrices from each individual
172 modelling method were aligned by covariate to create a single matrix containing an equal
173 number of bootstrap samples from each method. This combined matrix was used to calculate
174 overall covariate selection stabilities and Bootstrap P values as described above. Therefore
175 this method provided an estimate for selection stability and Bootstrap P values for all
176 covariates synthesised across all four methods of covariate selection.

177

178 2.4. Stability threshold calculation

179 A stability threshold above which each covariate was deemed to be important and selected in
180 a final model ('selection stability threshold (T^S)'), was calculated as follows. The concept
181 was to initially determine a distribution of covariate stability for each dataset on the basis that
182 there was *no casual relationship* between the outcome and explanatory variables and to
183 compare this to the stability of covariates in the actual data. The distribution of covariate
184 stability when no causal relationship existed we name here, the 'baseline stability
185 distribution'. When the stability of a covariate in the actual data had a specified probability
186 of exceeding values in the baseline stability distribution, it was deemed to be important and
187 selected in the final model.

188 For each model and dataset, the baseline stability distribution was determined by randomly
189 permuting the outcome variable to sever any relationship between outcome and explanatory
190 variables. Using the permuted outcome, a stability analysis was conducted as described in
191 Section 2.2 and the baseline stability of all covariates estimated. An empirical cumulative
192 distribution function (ECDF) of the baseline stability distribution was formulated and

193 stability thresholds were tested at different specified probability values of the ECDF. The
194 specified probability values were calculated as;

$$195 \qquad \qquad \qquad 1 - x/p \qquad \qquad \qquad (6)$$

196 where x was either equal to 0, 1, 3, 5, or 10 and p the number of covariates available in a
197 dataset. For example, when x = 1 and 500 potential covariates were available in the dataset,
198 the probability value of the ECDF at which the stability threshold would be set was;

$$199 \qquad \qquad \qquad 1 - (1/500) = 0.998$$

200 Therefore, x/p represented the probability that the stability of any covariate would exceed the
201 defined stability threshold.

202

203 To clarify the method, we list the steps used to calculate T^S ;

- 204 i. For a specified dataset and model, the outcome variable was randomly permuted to
205 remove the relationship between the outcome and explanatory variables.
- 206 ii. Covariate selection stability was conducted (as described in Section 2.2) using a set of
207 20 bootstrap samples.
- 208 iii. The distribution of covariate stability values for all ‘p’ covariates was defined using
209 an empirical cumulative distribution function (‘ecdf’ function, base R) and the (1-x/p)
210 probability value of the distribution calculated, to provide the stability threshold T^S .
211 The values tested for ‘x’ were 0, 1, 3, 5 and 10.
- 212 iv. Steps i. to iii. were repeated ten times and the mean and standard error for T^S
213 calculated across the ten repeats.
- 214 v. Covariate selection stability was conducted on the full dataset (i.e. without permuting
215 the outcome variable) using 200 bootstrap samples and the mean value of T^S , as
216 calculated in steps i – iv, used as the threshold above which variables were considered
217 to be ‘true’ and therefore selected in the final model.

218 Therefore, the stability threshold was set at a value at which there was a low probability (x/p)
219 of a covariate reaching this level of stability if it was not truly associated with the outcome.
220 Importantly, since location of the threshold is likely to be dependent upon the method of
221 modelling and specific dataset used, a new stability threshold was calculated for each dataset
222 and analysis undertaken.

223

224 2.5. Simulated datasets

225 The purpose of simulating data was to construct a set of realistic datasets, of varying size and
226 complexity, in which true underlying causal relationships were known. Therefore, model
227 performance could be evaluated objectively by comparing the accuracy with which each
228 method (including the stability threshold method) correctly selected the true causal variables.
229 For each dataset, 8 covariates were used to calculate an outcome variable; we refer to these as
230 the ‘true’ covariates. Additional noise variables were generated for each dataset, randomly
231 and independent of the outcome, we refer to these as ‘false’ variables. In each dataset, true
232 and false variables were set to be correlated with each other to a different extent such that
233 each statistical method could be evaluated with different but realistic complexity added to the
234 data (Table 1).

235 In all datasets, an outcome variable, “y_out”, was calculated from the true explanatory
236 variables as follows;

237

$$238 \quad y_out = Intercept + \sum_{k=1}^r 2.5 x_k + \sum_{j=1}^q 2.5 x_j + e \quad (7)$$

239

240 where *Intercept* =1, x_k represented the k^{th} of r true covariates simulated with no specified
241 correlation with other covariates, x_j represented the j^{th} of q true covariates that were specified
242 to have correlations with a set of ten false explanatory covariates, and e was a random

243 variable that represented all other true but unknown effects that causally influenced the
244 outcome, y_{out} . e was drawn from a Normal distribution with mean = 0 and standard
245 deviation which varied between dataset to allow the amount of variability explained by the
246 true covariates to differ between datasets (Table 1).

247 All explanatory covariates were drawn from distributions with mean = 0, SD = 1; they
248 represented standardised variables. Covariates that were simulated to be correlated (x_j) were
249 drawn from a multivariate normal distribution (using the `mvrnorm` function in the MASS
250 package²⁸ in R). In this case, each variable drawn from a distribution with mean = 0, SD = 1
251 and with covariance matrix specified such that variables were generated with a correlation
252 between 0.3 – 0.7, depending on the dataset (Table 1).

253 The characteristics of all datasets are summarised in Table 1. For each dataset, conventional
254 linear regression was used to model solely the true covariates, to illustrate the underlying
255 relationship between the outcome and true variables and estimate the total amount of
256 variation explained. True variables were simulated in all datasets such that their partial
257 coefficients in these regression models were individually significant ($P < 0.05$); that is their
258 signal was sufficiently strong to be deemed ‘significant’ when modelled alone. The purpose
259 of datasets containing more highly correlated covariates was to evaluate the extent to which
260 true variables could be retrieved when a substantial quantity of random and correlated noise
261 was added to the data and to explore whether selection stability enhanced retrieval.

262 For four of the 22 generated datasets a binary outcome was simulated instead of a continuous
263 outcome so that the analytic methods could be evaluated using an alternative outcome
264 distribution. A similar approach was used as in equation (6) except that an inverse logit
265 function followed by a draw from a binomial distribution was used to alter the outcome to be
266 1 or 0. As for the datasets with normally distributed outcomes, correlations were included
267 between true and false explanatory covariates to provide an increasing complexity within the

268 data. The variation explained by the true variables was estimated using an adapted R^2 ²⁹;
269 further details of the binomial outcome datasets are provided in Table 1.

270

271 {Table 1 approx here}

272

273 2.6. Evaluation and comparison of model performance

274 Model performance was assessed based on the correct selection of true variables and correct
275 omission of the false variables, and comparisons were made between methods for all datasets.
276 Graphical analysis was used to visualise model outcomes; scatterplots were used to illustrate
277 results of the stability analysis by plotting the stability (%) against BPV for covariates
278 selected in the model, as previously reported¹⁹ but with the additional annotation of the
279 selection stability threshold.

280

281 3. Results

282 3.1. Error rates for full (non-bootstrapped) models

283 A summary of the performance of all full (non-stability threshold) modelling methods, for
284 each dataset is provided in Table 2. Across all datasets, model selection using sAIC, elastic
285 net regression and MCP included a higher number of incorrect (false positive and false
286 negative) variables than selection using mBIC. The incorrect covariates selected using sAIC,
287 enet and MCP were mostly false positives indicating these models generally over fit the data.
288 Whilst the final models produced using mBIC were more sparse than for other methods
289 (fewer variables selected in the final models), the incorrectly classified covariates were
290 mostly false negatives indicating a tendency towards under fitting with this method. In
291 general, for most methods, more covariates were incorrectly classified as the complexity of
292 correlation structure within the data increased (Table 2).

293

294 {Table 2 approx here please}

295

296 3.2. Error rates for models incorporating covariate selection stability

297 Results of the covariate selection stability analyses using stability thresholds calculated with
298 $x = 0$, are provided in Table 3. Results for all other thresholds, $x = 1, 3, 5$ and 10 , are
299 provided in Supplementary Materials (Tables S1-S6). The stability thresholds that produced
300 fewest incorrectly classified covariates (false positives and false negatives) across all models
301 and all datasets were those calculated using $x = 0$ or $x = 1$ but the threshold calculated with x
302 $= 0$ selected fewest false positive covariates and was therefore judged the optimal threshold to
303 use. As expected, with increasing values of x , the number of false positive variables selected
304 tended to increase and the number of false negatives decreased. The one exception to $x=0$
305 being the optimal threshold was the highly complex dataset IC2, in which the stability
306 threshold at $x=3$ produced very slightly better results than $x=0$ or $x=1$ (Table S6).
307 Importantly, for all modelling methods, use of bootstrap selection stability at a threshold
308 using $x=0$ resulted in substantially fewer covariates being incorrectly selected (Table 3)
309 compared to the conventional implementation of each technique (Table 2). For sAIC, enet
310 and MCP the improvement in accuracy of covariate selection was marked; the number of
311 false positive covariates was greatly reduced with no increase in false negatives. For mBIC
312 the improvement in performance was less marked but still evident. The combination method,
313 a synthesis of methods incorporating results from all four model types, produced the best
314 overall performance compared to other methods, with the lowest number of false positive
315 variables selected whilst maintaining a low number of false negative covariates.

316

317 {Table 3 approx here}

318

319 Across all model types, it was noticeable that as complexity of correlations within the
320 simulated data increased, accuracy of covariate selection slightly deteriorated. Importantly, in
321 the datasets that were simulated with no or relatively few underlying correlations (Datasets
322 A-C), with a selection stability threshold set at $x = 0$ or $x = 1$, all model methods achieved
323 near-perfect selection accuracy when stability selection was used (Table 3). This was in
324 contrast to the performance of the models without selection stability (Table 2) in which only
325 mBIC models performed well on non-correlated data. As the size and number of correlations
326 within the data increased, performance of all methods slightly deteriorated, however, model
327 performance remained much better with the use of covariate selection stability than for
328 conventional implementation of each model type.

329

330 3.3. Diagrammatic illustration of model performance

331 Covariate selection stability thresholds for all values of x are illustrated graphically for
332 datasets MS1, MS2 and LS4, in Figures 1 – 3 respectively. The graphs provide an illustration
333 of the positioning of the stability thresholds (T^S) relative to the known true and false
334 covariates, for different values of x . Use of 10 replicates of 20 sets of bootstrap samples was
335 sufficient to produce relatively small 95% confidence intervals for T^S although the number of
336 replicates can be increased to further reduce the confidence interval as required.

337

338 {Figure 1 here, to be reproduced in colour please}

339

340 *Figure 1. Graphical representation of model results from dataset MS1. Each scatterplot*
341 *illustrates results of the stability analysis for one of five statistical methods; stepwise selection*
342 *using Akaike information criterion (AIC), stepwise selection using a modified Bayesian*

343 *Information Criterion (mBIC), elastic net regression (enet), Minimax Convex Penalty regression*
344 *(MCP) and the four methods combined (combi). The y-axis defines the negative bootstrap P value*
345 *for each covariate (-BPV) and the x-axis represents covariate stability (%). The dashed lines on*
346 *each graph (from right to left) represent the stability threshold calculated at x=0, 1, 3, 5 and 10*
347 *respectively. The shaded grey area depicts the 95% confidence interval for the stability threshold*
348 *calculated at x=0.*

349

350 {Figure 2 here please, to be reproduced in colour please}

351

352 *Figure 2. Graphical representation of model results from dataset MS2. Each scatterplot*
353 *illustrates results of the stability analysis for one of five statistical methods; stepwise selection*
354 *using Akaike information criterion (AIC), stepwise selection using a modified Bayesian*
355 *Information Criterion (mBIC), elastic net regression (enet), Minimax Convex Penalty regression*
356 *(MCP) and the four methods combined (combi). The y-axis defines the negative bootstrap P value*
357 *for each covariate (-BPV) and the x-axis represents covariate stability (%). The dashed lines on*
358 *each graph (from right to left) represent the stability threshold calculated at x=0, 1, 3, 5 and 10*
359 *respectively. The shaded grey area depicts the 95% confidence interval for the stability threshold*
360 *calculated at x=0.*

361

362 {Figure 3 here, to be reproduced in colour please}

363

364 *Figure 3. Graphical representation of model results from dataset LS4. Each scatterplot illustrates*
365 *results of the stability analysis for one of five statistical methods; stepwise selection using Akaike*
366 *information criterion (AIC), stepwise selection using a modified Bayesian Information Criterion*
367 *(mBIC), elastic net regression (enet), Minimax Convex Penalty regression (MCP) and the four*
368 *methods combined (combi). The y-axis defines the negative bootstrap P value for each covariate*

369 *(-BPV) and the x-axis represents covariate stability (%). The dashed lines on each graph (from*
370 *right to left) represent the stability threshold calculated at $x=0, 1, 3, 5$ and 10 respectively. The*
371 *shaded grey area depicts the 95% confidence interval for the stability threshold calculated at*
372 *$x=0$.*

373

374 4. Discussion

375 Our results using simulated datasets with known relationships highlight the value of selection
376 stability as a means to identify the correct explanatory variables when the number of
377 covariates (p) is relatively large with respect to the number of observations (n). Whilst the
378 concept of selection stability is not new^{9,27}, its use within epidemiological fields has been
379 relatively limited. It is generally considered that selection stability is a necessary and
380 beneficial approach to enhance inferential modelling with high dimensional data^{11,30} and that
381 without selection stability results of different selection methods vary dramatically leading to
382 confusion in inference^{26,31}.

383 An important consideration, however, when employing selection stability is to identify a
384 stability threshold above which covariates are deemed ‘important’ or ‘significant’ and
385 selected in a final model; this was the primary purpose of the current study. Whilst it is
386 known that a higher threshold for selection stability will result in the selection of fewer false
387 positive covariates⁹, a generalisable approach to identify an optimal threshold suitable for
388 multiple model types has not been reported. A threshold proposed for one model type, lasso
389 regression⁹ has been shown to be too conservative with too few true variables being selected
390^{13,32}. A further problem with the threshold proposed by Meinshausen and Bühlmann⁹ is that it
391 cannot be defined at a value $<50\%$ which, as we have shown in this research, is certainly
392 necessary with some types of models and data (Figures 1-3).

393 In this study we defined and evaluated a straightforward approach to identify a covariate
394 selection stability threshold that can be readily used with different datasets and model types.
395 The intuitive approach of using an expected or ‘baseline’ stability distribution from a model
396 with no underlying relationships and using the same data to compare this to a distribution
397 from a model with causal relationships, worked well across all five model types and all 22
398 complex datasets. Even for the highly correlated datasets, error rates were relatively low and,
399 importantly, were markedly superior to the equivalent conventional (non-stability selection)
400 models. Therefore, for safety of inference, when using relatively wide, high dimensional
401 datasets, covariate selection stability with the defined threshold will minimise identification
402 of false positive and false negative variables compared to conventional modelling techniques.
403 Furthermore, although highly correlated data are generally recognised to distort model
404 parameters³³, our results indicate that selection stability with an appropriate choice of
405 threshold provides a method to mitigate this issue.

406 The concept of permuting variables to enhance statistical inference is not new. It has been
407 used in machine learning to facilitate estimation of variable importance^{34,35} and as a
408 replacement for cross validation to identify covariates in elastic net regression using non-
409 sparse genetic data¹³. In the setting we describe, permuting the outcome allowed
410 identification of an optimal covariate stability threshold to identify true variables for a variety
411 of modelling methods and types of data. Whilst using covariate selection stability was more
412 computationally expensive than conventional approaches, this was not prohibitive. For
413 example, for a set of 10 repeats of 20 bootstrap samples (to define the baseline stability
414 threshold) followed by 200 bootstrap samples (to select variables above the baseline
415 threshold) using a dataset with 500 observations and 1000 covariates, the total time taken for
416 computation of the mBIC method was 3 minutes and 16 seconds using a 15 core i9 processor
417 with 32GB RAM. The equivalent time for the elastic net model was 28 minutes and 4

418 seconds. In terms of the substantial improvements in accuracy of covariate selection, the
419 additional time required would appear to be very worthwhile.
420 Also worthy of note in our results was the excellent performance of the combination
421 modelling method; this is in agreement with previous research ¹⁹. For all datasets this method
422 provided the most accurate variable selection and this is likely to be due to the principle of
423 triangulation. Triangulation is based on the concept that reliability of results is enhanced
424 through integration of different approaches, particularly when each approach has a different,
425 unrelated source of bias ³⁶. For this reason, it has been suggested that triangulation should be
426 more widely adopted to ensure robustness and reproducibility of scientific results ³⁷ although
427 it has rarely been employed in animal health research.

428

429 4.1. Study limitations

430 There were several limitations of this study worthy of consideration. Although substantial
431 complexity was incorporated into the simulated datasets, no clustering of the outcome
432 variable was considered. Therefore, extension of this stability threshold method to mixed
433 effect models, although likely to be valid, requires additional investigation. Similarly, non-
434 linearities within predictor variables were not considered in the simulated data and this also
435 warrants additional exploration in terms of selection stability thresholds. Four methods of
436 modelling were evaluated, and a combination of all four together, which suggests our
437 proposed method of selection stability may be generalisable, however, further research using
438 additional model types would be worthwhile. In addition, although one group of simulated
439 dataset incorporated discrete (binary) outcome variables (Datasets BIN1-4), further research
440 would be valuable to investigate the value of these methods with outcome distributions of
441 varying types. Finally, it should be noted that underlying true solutions in our simulated data
442 were relatively sparse; eight covariates were chosen to have true causal influences. This

443 aligns with the ‘bet on sparsity principle’⁴, which dictates that for causal variables to be
444 identifiable (in any model), a relatively small number of predictors have to be responsible for
445 most of the effect on the outcome of interest; this is why such data structures were chosen in
446 this study.

447

448 4.2. Conclusions

449 In this research we have proposed a new method to conduct covariate selection for use with
450 relatively wide data. The method is based on a new approach to covariate selection stability
451 that incorporates a stability threshold to define which covariates should be included in a final
452 model. Results indicate the approach offers substantial reductions in covariate selection error
453 rates compared with conventional model selection methods.

454

455 **Acknowledgements**

456 The author Eliana Lima is employed with the European Food Safety Authority (EFSA) in the
457 Unit ALPHA that provides scientific and administrative support to EFSA's Scientific
458 Activities in the area of animal health and welfare. However, the present article is published
459 under the sole responsibility of the author Eliana Lima and may not be considered as an
460 EFSA scientific output. The positions and opinions presented in this article are those of the
461 author alone and are not intended to represent any views/any official position or scientific
462 work of EFSA. To know about the views or scientific outputs of EFSA, please consult its
463 website under <http://www.efsa.europa.eu>.

464 **References**

- 465 1. Wasserman, L. & Roeder, K. High Dimensional Variable Selection. *Ann. Stat.* **1**,
466 2178–2201 (2009).
- 467 2. Sirimongkolkasem, T. & Drikvandi, R. On Regularisation Methods for Analysis of
468 High Dimensional Data. *Ann. Data Sci.* **6**, 737–763 (2019).
- 469 3. Liu, J. Y., Zhong, W. & Li, R. Z. A selective overview of feature screening for
470 ultrahigh-dimensional data. *Sci. China Math.* **58**, 2033–2054 (2015).
- 471 4. Statistical Learning with Sparsity: the Lasso and Generalizations. Available at:
472 <https://web.stanford.edu/~hastie/StatLearnSparsity/>. (Accessed: 2nd March 2020)
- 473 5. Bogdan, M., Ghosh, J. K. & Zak-Szatkowska, M. Selecting explanatory variables with
474 the modified version of the bayesian information criterion. in *Quality and Reliability*
475 *Engineering International* **24**, 627–641 (2008).
- 476 6. Fan, J. & Peng, H. Nonconcave penalized likelihood with a diverging number of
477 parameters. *Ann. Stat.* **32**, 928–961 (2004).
- 478 7. Tibshirani, R. Regression Shrinkage and Selection Via the Lasso. *J. R. Stat. Soc. Ser. B*
479 **58**, 267–288 (1996).
- 480 8. Zou, H. & Hastie, T. Regularization and variable selection via the elastic net. *J. R.*
481 *Stat. Soc. Ser. B (Statistical Methodol.* **67**, 301–320 (2005).
- 482 9. Meinshausen, N. & Bühlmann, P. Stability selection. *J. R. Stat. Soc. Ser. B (Statistical*
483 *Methodol.* **72**, 417–473 (2010).
- 484 10. Baldassarre, L., Pontil, M. & Mourão-miranda, J. Sparsity Is Better with Stability :
485 Combining Accuracy and Stability for Model Selection in Brain Decoding. *Front.*
486 *Neurosci.* **11**, 62 (2017).
- 487 11. Heinze, G., Wallisch, C. & Dunkler, D. Variable selection – A review and
488 recommendations for the practicing statistician. *Biometrical Journal* **60**, 431–449

- 489 (2018).
- 490 12. Alexander, D. H. & Lange, K. Stability selection for genome-wide association. *Genet.*
491 *Epidemiol.* **35**, 722–728 (2011).
- 492 13. Kim, K., Koo, J. & Sun, H. An empirical threshold of selection probability for analysis
493 of high-dimensional correlated data. *J. Stat. Comput. Simul.* **90**, 1606–1617 (2020).
- 494 14. Lima, E. *et al.* Use of bootstrapped, regularised regression to identify factors
495 associated with lamb-derived revenue on commercial sheep farms. *Prev. Vet. Med.*
496 **174**, 104851 (2020).
- 497 15. Hyde, R. M., Green, M. J., Hudson, C. & Down, P. M. Factors associated with daily
498 weight gain in preweaned calves on dairy farms. *Prev. Vet. Med.* **190**, 105320 (2021).
- 499 16. Spooner, A. *et al.* A comparison of machine learning methods for survival analysis of
500 high-dimensional clinical data for dementia prediction. *Sci. Rep.* **10**, 20410 (2020).
- 501 17. Tibshirani, R. *et al.* *An Introduction to Statistical Learning with*
502 *Applications in R (older version). Springer Texts in Statistics* (New York : Springer,
503 [2013] ©2013, 2013).
- 504 18. Zhang, C. H. Nearly unbiased variable selection under minimax concave penalty. *Ann.*
505 *Stat.* **38**, 894–942 (2010).
- 506 19. Lima, E., Hyde, R. & Green, M. Model selection for inferential models with high
507 dimensional data: synthesis and graphical representation of multiple techniques. *Sci.*
508 *Rep.* **11**, 412 (2021).
- 509 20. Piotr Szulc. bigstep: Stepwise Selection for Large Data Sets. R package. (2019).
- 510 21. R Core Team. R: A language and environment for statistical computing (version
511 1.1.463). (2018).
- 512 22. Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear
513 Models via Coordinate Descent. *J. Stat. Softw.* **33**, 1–22 (2010).

- 514 23. Kuhn, M. *et al.* caret: Classification and Regression Training. (2019).
- 515 24. Breheny, P. & Huang, J. Coordinate descent algorithms for nonconvex penalized
516 regression, with applications to biological feature selection. *Ann. Appl. Stat.* **5**, 232–
517 253 (2011).
- 518 25. Breiman, L. Bagging Predictors. *Mach. Learn.* **24**, 123–140 (1996).
- 519 26. Lima, E., Davies, P., Kaler, J., Lovatt, F. & Green, M. Variable selection for
520 inferential models with relatively high-dimensional data: Between method
521 heterogeneity and covariate stability as adjuncts to robust selection. *Sci. Rep.* **10**, 1–11
522 (2020).
- 523 27. Sauerbrei, W. The Use of Resampling Methods to Simplify Regression Models in
524 Medical Statistics. *J. R. Stat. Soc. Ser. C (Applied Stat.)* **48**, 313–329 (1999).
- 525 28. Venables, W. N. & Ripley, B. D. *Modern applied statistics with S. Statistics and
526 computing* **45**, (Springer-Verlag New York, 2002).
- 527 29. Nagelkerke, N. J. D. A note on a general definition of the coefficient of determination.
528 *Biometrika* **78**, 691–692 (1991).
- 529 30. Wang, F., Mukherjee, S., Richardson, S. & Hill, S. M. High-dimensional regression in
530 practice: an empirical study of finite-sample prediction, variable selection and ranking.
531 *Stat. Comput.* **30**, 697–719 (2020).
- 532 31. Pfeiffer, R. M., Redd, A. & Carroll, R. J. On the impact of model selection on
533 predictor identification and parameter inference. *Comput. Stat.* **32**, 667–690 (2017).
- 534 32. Hofner, B., Boccuto, L. & Göker, M. Controlling false discoveries in high-dimensional
535 situations: Boosting with stability selection. *BMC Bioinformatics* **16**, 144 (2015).
- 536 33. Dormann, C. F. *et al.* Collinearity: a review of methods to deal with it and a simulation
537 study evaluating their performance. *Ecography (Cop.)*. **36**, 27–46 (2013).
- 538 34. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).

- 539 35. Kursa, M. B., Jankowski, A. & Rudnicki, W. R. Boruta - A system for feature
540 selection. *Fundam. Informaticae* **101**, 271–285 (2010).
- 541 36. Lawlor, D. A., Tilling, K. & Smith, G. D. Triangulation in aetiological epidemiology.
542 *Int. J. Epidemiol.* **45**, 1866–1886 (2016).
- 543 37. Munafò, M. R. & Davey Smith, G. Robust research needs many lines of evidence.
544 *Nature* **553**, 399–401 (2018).
- 545
- 546

547 Author contributions statement

548 EL, RH and MG contributed to development of statistical techniques and statistical analysis.

549 MG, RH and EL all contributed to the writing of the final manuscript.

550 Table 1. Description of the simulated datasets used to evaluate five modelling techniques
 551 incorporating covariate selection stability.

552
 553

Dataset name	p	True covariates	n	Description of correlation structure	R ²
HS: CONTINUOUS NORMALLY DISTRIBUTED OUTCOME; HIGH SIGNAL					
HS1	500	8	500	all true and false variables uncorrelated	0.70
HS2	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.3	0.71
HS3	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.5	0.69
HS4	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.7	0.71
MS: CONTINUOUS NORMALLY DISTRIBUTED OUTCOME ; MEDIUM SIGNAL					
MS1	500	8	500	all true and false variables uncorrelated	0.51
MS2	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.3	0.46
MS3	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.5	0.45
MS4	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.7	0.45
LS: CONTINUOUS NORMALLY DISTRIBUTED OUTCOME; LOW SIGNAL					
LS1	500	8	500	all true and false variables uncorrelated	0.29
LS2	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.3	0.26
LS3	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.5	0.28
LS4	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.7	0.29
WID: CONTINUOUS NORMALLY DISTRIBUTED OUTCOME; WIDER DATASET					
WID1	1000	8	500	all true and false variables uncorrelated	0.51
WID2	1000	8	500	4 true covariates correlated individually with 10 false covariates at 0.3	0.46
WID3	1000	8	500	4 true covariates correlated individually with 10 false covariates at 0.5	0.45
WID4	1000	8	500	4 true covariates correlated individually with 10 false covariates at 0.7	0.45
BIN: BINARY OUTCOME MEDIUM SIGNAL					
BIN1	500	8	500	all true and false variables uncorrelated	0.48
BIN2	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.3	0.48
BIN3	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.5	0.47
BIN4	500	8	500	4 true covariates correlated individually with 10 false covariates at 0.7	0.47
IC: CONTINUOUS NORMALLY DISTRIBUTED OUTCOME; INCREASED COMPLEXITY					
IC1	500	8	500	4 true covariates correlated with each other and 20 false covariates at 0.9	0.56
IC2	500	8	500	4 true covariates correlated with each other and 92 false covariates at 0.7	0.42

554
 555
 556
 557
 558
 559

Key: p – total number of covariates offered to the model, True covariates – number of covariates that were used to directly calculate the outcome and therefore deemed causal, n – number of observations (rows) of data, R² – amount of total variation in the outcome variable explained by the true covariates (for binary outcome data, an adapted R² was used (Nagelkerke, 1991).

560 Table 2. Results of full (non-bootstrapped) models conducted on 22 simulated datasets (Table
 561 1) to illustrate the number of false positive and false negative covariates selected in final models
 562 for each method.

563
 564

Dataset name	Model type							
	sAIC		mBIC		enet		MCP	
	FP	FN	FP	FN	FP	FN	FP	FN
HIGH SIGNAL								
HS1	11	1	0	1	22	0	1	1
HS2	15	1	0	1	31	0	2	1
HS3	15	1	0	1	31	0	3	1
HS4	17	1	0	1	35	0	1	1
MEDIUM SIGNAL								
MS1	14	1	0	1	9	0	7	1
MS2	13	1	0	1	11	0	6	1
MS3	25	1	0	1	27	0	15	1
MS4	25	1	0	1	27	0	4	1
LOW SIGNAL								
LS1	15	1	0	3	38	0	17	1
LS2	23	1	0	3	48	0	28	1
LS3	19	1	0	2	52	0	14	1
LS4	31	1	1	4	57	0	15	2
WIDER DATASET								
WID1	24	1	0	1	12	0	6	1
WID2	24	1	0	1	19	0	1	1
WID3	39	1	0	1	66	0	8	1
WID4	46	1	0	1	45	0	8	1
BINARY OUTCOME								
BIN1	17	1	0	1	20	0	5	1
BIN2	17	1	0	1	33	0	8	1
BIN3	28	1	0	1	94	0	9	1
BIN4	20	2	1	2	43	0	8	2
INCREASED COMPLEXITY								
IC1	15	1	0	2	30	0	24	2
IC2	25	2	0	3	51	0	14	2

565

566

567 Key; sAIC – covariate selection based on stepwise Akaike information criterion, mBIC - covariate selection
 568 based on a stepwise on a modified Bayesian Information Criterion, enet – covariate section based on elastic net
 569 regression, MCP – covariate selection based on minimax convex penalty regression, FP – number of false
 570 positive covariates selected, FN – number of false negative variables selected

571

572

573

574 Table 3. Results of models incorporating covariate selection stability (with selection threshold
 575 set at $\alpha = 0$), conducted on 22 simulated datasets (Table 1), to illustrate the number of false
 576 positive and false negative covariates selected in final models for each method.

577
 578

Dataset name	Model type									
	sAIC		mBIC		enet		MCP		Combi	
	FP	FN	FP	FN	FP	FN	FP	FN	FP	FN
HIGH SIGNAL										
HS1	0	0	0	0	0	0	0	0	0	0
HS2	0	0	0	0	0	0	0	0	0	0
HS3	0	0	0	0	0	0	0	0	0	0
HS4	0	0	1	0	0	0	0	0	0	0
MEDIUM SIGNAL										
MS1	0	0	0	0	0	0	0	0	0	0
MS2	0	0	0	0	0	0	0	0	0	0
MS3	0	0	0	0	0	0	0	0	0	0
MS4	1	0	0	0	0	0	0	0	0	0
LOW SIGNAL										
LS1	0	2	0	1	0	2	0	1	0	0
LS2	0	2	1	0	0	2	0	2	0	0
LS3	0	3	2	0	0	2	0	1	1	0
LS4	1	1	1	1	0	2	0	1	0	1
WIDER DATASET										
WID1	0	0	0	0	0	0	0	0	0	0
WID2	0	0	0	0	0	0	0	0	0	0
WID3	0	0	0	0	0	0	0	0	0	0
WID4	0	0	0	0	0	0	0	0	0	0
BINARY OUTCOME										
BIN1	0	0	2	0	0	0	0	0	0	0
BIN2	0	1	1	0	0	0	0	0	0	0
BIN3	1	0	1	0	0	0	0	0	1	0
BIN4	0	1	1	1	0	1	0	1	0	1
INCREASED COMPLEXITY										
IC1	0	4	1	0	0	2	0	3	0	3
IC2	0	3	0	2	0	3	0	4	0	3

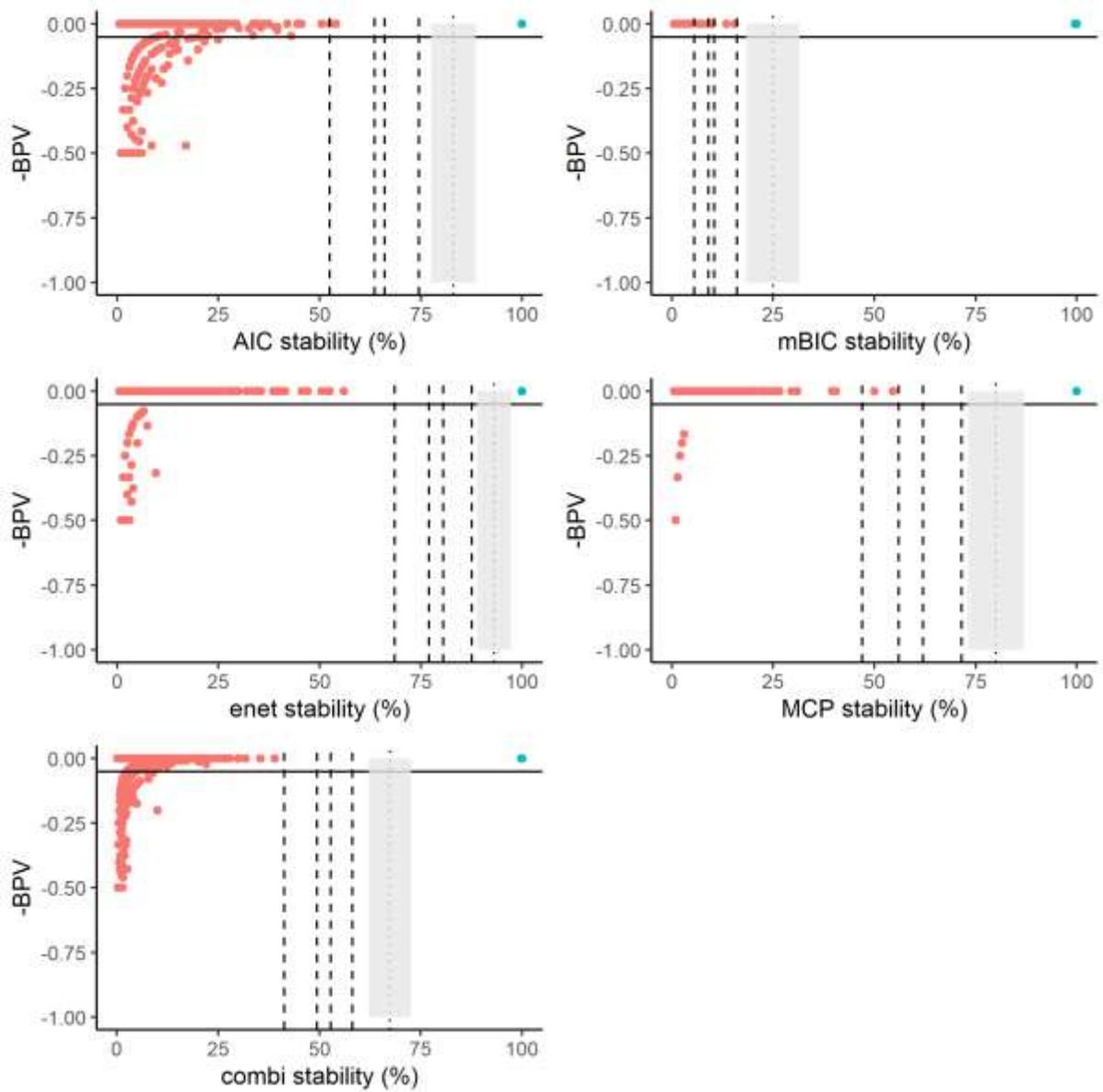
579
 580

581 Key; sAIC – covariate selection based on stepwise Akaike information criterion, mBIC - covariate selection
 582 based on a stepwise on a modified Bayesian Information Criterion, enet – covariate section based on elastic net
 583 regression, MCP – covariate selection based on minimax convex penalty regression, Combi – covariate
 584 selection based on a synthesised combination of sAIC, mBIC, enet and MCP, FP – number of false positive
 585 covariates selected, FN – number of false negative variables selected

586

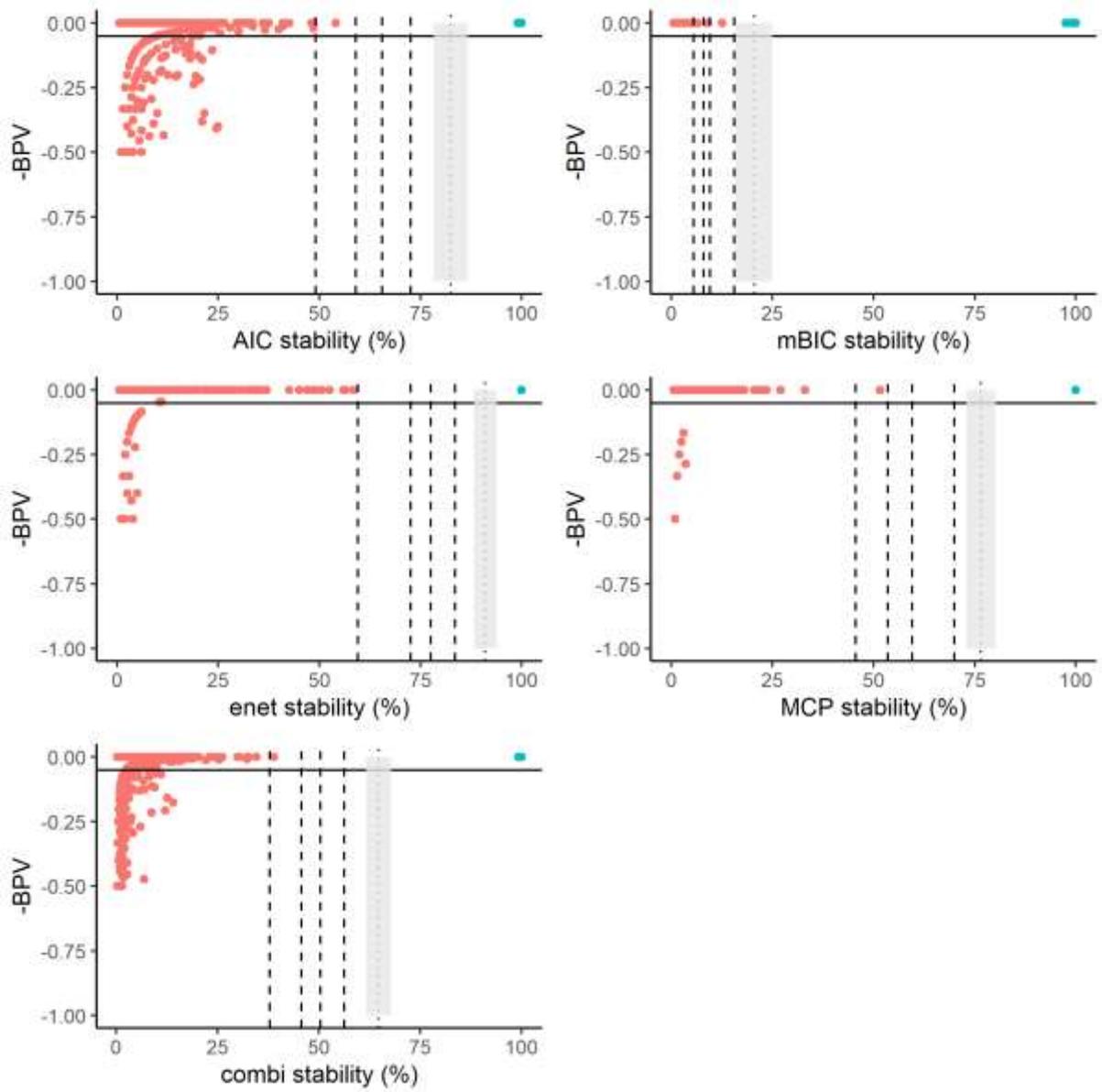
587
588
589
590

Figure 1.



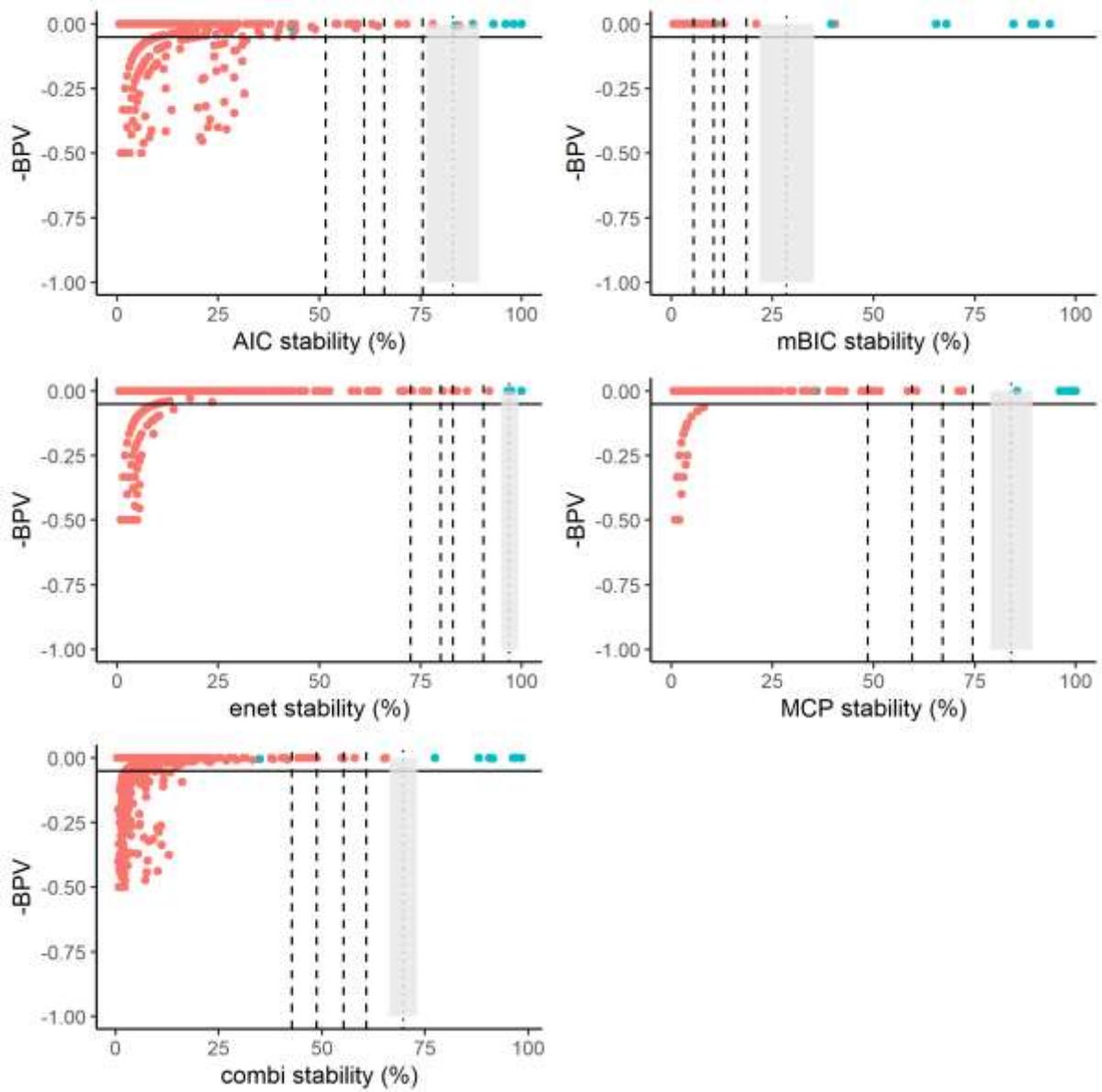
591
592
593
594
595
596
597
598
599
600
601
602
603

604 Figure 2.
605
606



607
608
609
610
611
612
613
614

615 Figure 3.
616



617
618
619

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementarymaterials.docx](#)