1     **Towards the systematic simplification of mechanistic models**

2

3     [1]G. M. Cox, [1]J. M. Gibbons, [2]A. T. A. Wood, [1]J. Craigon, [1]S. J. Ramsden

4     and [1]N. M. J. Crout.

5

6     [1]School of Biosciences, University of Nottingham, Sutton Bonington, LE12 5RD.

7     [2]School of Mathematical Sciences, University of Nottingham, University Park,

8     Nottingham, NG7 2RD.

9

10    **Abstract**

11    Mechanistic models used for prediction should be parsimonious, as models

12    which are over-parameterised may have poor predictive performance.

13    Determining whether a model is parsimonious requires comparisons with

14    alternative model formulations with differing levels of complexity.

15    However, creating alternative formulations for large mechanistic models is

16    often problematic, and usually time-consuming. Consequently, few are

17    ever investigated. In this paper, we present an approach which rapidly

18    generates reduced model formulations by replacing a model's variables

19    with constants. These reduced alternatives can be compared to the

20    original model, using data based model selection criteria, to assist in the

21    identification of potentially unnecessary model complexity, and thereby

22    inform reformulation of the model. To illustrate the approach, we present

23    its application to a published radiocaesium plant-uptake model, which

24    predicts uptake on the basis of soil characteristics (e.g. pH, organic matter

25    content, clay content). A total of 1024 reduced model formulations were

26  generated, and ranked according to five model selection criteria: Residual

27  Sum of Squares (RSS), $AIC_c$, BIC, MDL  and ICOMP. The lowest scores for

28  RSS and $AIC_c$ occurred for the same reduced model in which pH

29  dependent model components were replaced. The lowest scores for BIC,

30  MDL and ICOMP occurred for a further reduced model in which model

31  components related to the distinction between adsorption on clay and

32  organic surfaces were replaced. Both these reduced models had a lower

33  RSS for the parameterisation dataset than the original model. As a test of

34  their predictive performance, the original model and the two reduced

35  models outlined above were used to predict an independent dataset. The

36  reduced models have lower prediction sums of squares than the original

37  model, suggesting that the latter may be overfitted. The approach

38  presented has the potential to inform model development by rapidly

39  creating a class of alternative model formulations, which can be

40  compared.

41

42  **Introduction**

43  Mechanistic, or process based, models are generally highly structured and

44  have inter-related components whose mathematical specification is

45  informed by scientific knowledge of relevant processes. Models of this type

46  are widely used. Mechanistic models are usually developed using expert

47  knowledge of the processes involved in the system under consideration.

48  This development may include the amalgamation of previously established

49  relationships (e.g. Gibbons et *al*., 2005), the development of new

50  relationships (e.g. Crout *et al*., 1998), or, more commonly, a combination

51  of both (e.g. Jamieson *et al*., 1998). If an appropriate dataset is available,

52  the model parameters may be chosen to achieve the best "fit", in which

53  case the model may be described as being semi-mechanistic. If parameter

54  values are chosen using a numerical procedure (e.g. least squares), we

55  term this "formal parameterisation". Often, if the goodness-of-fit (GOF) is

56  considered inadequate, the model may be modified by the addition of new

57  parameters or relationships. Throughout this development process,

58  judgements (which are often implicit) are made about the appropriate

59  level of complexity in the model.  However, it is well known that a model's

60  fit to a particular dataset can always be improved by the addition of new

61  parameters, and that this may lead to over-fitting and poor predictive

62  performance when the model is applied to a new situation (e.g. Myung

63  and Pitt, 2002). To avoid these difficulties model developers may adhere

64  to the parsimony principle, which states that "models should be as simple

65  as possible, but no simpler". Unfortunately, determining the point of

66  optimal model simplicity is often difficult in practice, as this requires the

67  generation and comparison of alternative model formulations. Generating

68  alternative formulations of large mechanistic or semi-mechanistic models

69  may not be straightforward, and can be very time-consuming.

70  Consequently, although there are often many plausible representations of

71  a given system, simpler alternatives are rarely assessed. This is in sharp

72  contrast with, for example, linear statistical models for which coefficients

73  can be readily set to zero to investigate reduced models.

74  One approach to creating a set of alternative models is "model

75  generation". For example, Atanasova *et al*. (2006) describe an automated

76  modelling tool where experts define a "knowledge library" containing

77 context free grammar statements that characterise the general processes

78 involved in the system under study. Different models are generated by

79 combining the various expressions specified for each general process. The

80 models are then parameterised by the fitting of constants, and the best

81 performing models identified.

82 A limitation of such approaches is that for complex systems, where there

83 may exist many alternative explanations of the underlying processes, the

84 number of possible models can be very large, rendering parameterisation

85 of the candidate models infeasible.

86 More recently, Asgharbeygi *et al*. (2006) have developed an algorithm

87 which generates a set of alternative models based on an initial model, i.e.

88 "model revision". Users specify which parts of the initial model are "fixed"

89 and which parts can be removed or have their parameters changed,

90 reflecting the areas of uncertainty within the model. The algorithm

91 generates all models that are consistent with the constraints specified,

92 and each model structure is parameterised using observed data. The

93 method we describe here is similar, although simpler, and we are focussed

94 upon the systematic removal of variables from a model, rather than the

95 insertion or alteration of processes.

96 We illustrate our approach through its application to a published model,

97 and discuss the results both in the context of the example model and

98 more general application.

99

**100**   **Approach**

101   Before describing the approach, we define some terminology. Constant

102   values within a model are *parameters*. For the purposes of the model

103   development, they may be fixed, in which case their value is set before

104   the model was developed, or they may be adjustable in which case their

105   value is estimated as part of the model development process, usually

106   through the use of data. *Input variables* are values obtained directly from

107   data, and are independent of a model's calculations. *Model variables* are

108   internal quantities calculated using an assumed relationship expressed in

109   terms of the model's parameters, input variables and other model

110   variables. The definition of model variables is partially subjective because

111   intermediate steps in a model calculation could be defined as individual

112   model variables, or combined into a larger relationship as a single model

113   variable. Such choices will often depend upon the requirements of specific

114   computer implementation. However, for our purposes, we shall regard

115   each model variable as having a specific mechanistic interpretation. This is

116   illustrated later in the example application. Throughout we use M to

117   denote the number of model variables, p to denote the number of

118   parameters and n to denote the number of data.

119   Traditional statistical approaches to model selection have focussed on the

120   number of adjustable parameters as a measure of model complexity

121   (either explicitly or implicitly). Here we are also considering the number of

122   model variables and inputs as a further measure of model complexity in

123   order to reflect the structured and inter-related nature of typical

124   mechanistic models. This distinction is further illustrated with reference to

125   the example we present later.

126   The approach investigated involves the systematic replacement of model

127   variables by constant values to produce a class of reduced models. The

128   performance of these reduced models can then be compared using various

129   criteria to assist the identification of model variables whose inclusion are

130   not justified by the data, and which may, therefore, be unnecessarily

131   increasing the complexity of the model. The procedure is not intended to

132   generate the *best* model, rather, it is hoped that it may be used as an

133   iterative diagnostic to inform model development.

134   Consider a model comprised of M model variables, $V_i$, each of which is

135   defined by a relationship in terms of parameters, input variables or other

136   model variables. If all of the possible combinations of variable

137   replacements, $R_i$, are considered (i.e. an exhaustive search), $2^M$ simplified

138   models will be generated and require assessment. If the model considered

139   contains parameters which have been estimated using data then it may be

140   appropriate to re-estimate these values for each reduced model.

141

142   <u>Choice of replacement value</u>

143   An important question when simplifying mechanistic models by replacing

144   model variables with constants is: how should the replacement values be

145   selected? In principle, our objectives could be met by setting $R_i$ to

146   arbitrary values. However, the $R_i$ need to be chosen in such a way that the

147   rest of the model calculations can proceed successfully. A feature of many

148   mechanistic models is the high degree of inter-connection between model

149   variables, where one variable may depend upon another and so on.

150   Consequently, an inappropriate choice of $R_i$ may lead to poor model

151  performance and/or numerical problems (e.g. if the value of the

152  replacement constant results in taking the logarithm of a negative

153  number). For this reason the standard approach for linear models, in

154  which coefficients are set to zero, is not appropriate. One practical method

155  is to set $R_i$ equal to the mean value $V_i$ attains over the course of a

156  simulation in which there are no replacements (i.e. using the original

157  model). The rationale for this method is that the replacement value is

158  broadly appropriate, and our comparison between models becomes a test

159  of whether the variation of a model variable about its mean is worth

160  including in the model.

161  An obvious temptation here would be to select values for the $R_i$, via formal

162  parameterisation, which maximised the likelihood function. However,

163  whilst this would improve the fit of the reduced models, it would

164  effectively be introducing new adjustable parameters and consequently

165  increase model complexity. This would conflict with our objective of

166  identifying parsimonious models.

167  A further problem with using fitted replacement constants is that they

168  may make interpretation of the results more difficult if the optimised

169  values obtained are not mechanistically feasible. This can be avoided if the

170  parameters' values are constrained in some way, although, care must

171  taken when defining parameter boundaries, as limits which are too

172  restrictive may affect the predictive performance of any reduced models

173  generated. A further limitation to this approach is that it is

174  computationally more intensive than simply using mean values, due to the

175  fitting of the replacements. This may be significant when performing

176    exhaustive searches with many replacement candidates, especially for

177    large models.

178

179    <u>Comparing Model Performance</u>

180    The ideal measure of a model's predictive performance is how well it can

181    predict observed values of interest for a new situation. When a suitable

182    dataset, which has not been used for model development, is available its

183    predictive performance can be assessed by a measure such as the

184    prediction residual sum of squares (PSS), defined as the sum of squared

185    differences between the observed and predicted values.

186    If independent data are not available, an alternative approach is to rely on

187    RSS (or other GOF statistics) derived using the data employed during

188    model development. However, as discussed earlier, this does not take into

189    account the possibility that the model is over-fitted. In these cases model

190    selection criteria are a useful alternative, although it should be noted that

191    they are only applicable if the model has been formally parameterised.

192    Several model selection criteria have been developed in the fields of

193    information science and statistics, some of which are summarised in Table

194    1. Each comprises a term based on the model's GOF and a term which

195    estimates the influence of the model's complexity on its predictive

196    capability.

197    The models we consider are all of the following general form:

198    $$y_j = f(I_j, \theta) + \in_j, \quad j = 1,...,n,$$

199    where $n$ is the sample size, $y_j$ is the response for observational unit $j$, $I_j$ is

200    the corresponding vector of values of the input variables, $\theta$ is the

201    parameter vector for the model under consideration, $f$ is a known function

202    of $I_j$ and $\theta$, and $\in_1, ..., \in_n$ are independent random error terms which are

203    normally distributed with mean zero and variance $\sigma^2$. Each model

204    determines an $f$. For the models under consideration, $f$ is too complicated

205    to specify explicitly here; an idea of the structure of a typical $f$ is given by

206    Figure 1. In practice, each $f$ is specified through a computer program.

207    The log-likelihood for a model is given by

208    
$$l(\theta, \sigma^2) = -\frac{n}{2}\ln(\sigma^2) - \frac{1}{2\sigma^2}\sum_{j=1}^{n}(y_j - f(I_j, \theta))^2$$

209    The maximised log-likelihood is given by

210    
$$\ln(ML) = l(\hat{\theta}, \hat{\sigma}^2),$$

211    where $\hat{\theta}$ is found by numerically by using the Marquardt parameterisation

212    procedure (Press *et al.*, 1989) which minimises the residual sum of

213    squares

214    
$$\sum_{j=1}^{n}(y_i - f(I_j, \theta))^2,$$

215    and $\hat{\sigma}^2$ is set equal to the minimised mean residual sum of squares, i.e.

216    
$$\hat{\sigma}^2 = \frac{1}{n}\sum_{j=1}^{n}(y_i - f(I_j, \hat{\theta}))^2.$$

217    The principal difference between the model selection criteria is the

218    approach used to estimate model complexity. In AIC (Akaike's Information

219    Criterion), the complexity term is simply twice the number of adjustable

220    parameters in the model. However, where sample sizes are small,

221    Burnham and Anderson (2002) recommend using $AIC_c$, a corrected

222    version of AIC, when n/p<40. In BIC (Bayesian Information Criterion), the

223    number of data points used to calculate the maximum likelihood is

224  introduced, and consequently BIC penalises parameters more than AIC

225  when n>8. However, complexity may not be related simply to the *number*

226  of parameters in a model, but also the model's functional form. The MDL

227  (Minimum Description Length) and ICOMP (Information Complexity)

228  criteria attempt to take this into account through the Hessian matrix

229  (which is the matrix of second derivatives, with respect to $\theta$ and $\sigma^2$, of

230  the log-likelihood $l(\theta, \sigma^2)$, evaluated at $\theta = \hat{\theta}$ and $\sigma^2 = \hat{\sigma}^2$) and the

231  asymptotic covariance matrix of the parameter estimates respectively.

232  These matrices are estimated during the Marquardt parameterisation

233  procedure.

234  In the context of the model selection criteria, only adjustable parameters

235  that are estimated using data are considered when determining the level

236  of model complexity. However, determining the number of parameters to

237  be included within the criteria may not be straightforward, as frequently

238  some "fixed" parameters (which are not included in formal

239  parameterisation procedures) are "tweaked" (i.e. adjusted manually by

240  model developers) during model development to obtain a better fit, which

241  amounts to *ad hoc* parameterisation. If that is the case, those parameters

242  should be considered by the selection criteria.

243  Finally, it should be noted that the derivations of these selection criteria

244  include a series of simplifying assumptions, which may not be satisfied in

245  all cases. Consequently, some caution is required in the application of

246  these measures. See Burnham and Andersen (2002) and Raftery (1995)

247  for relevant discussion.

248

## 249    Example Application

250    <u>Model description</u>

251    The model developed by Absalom *et al*. (2001) predicts the plant uptake

252    of radiocaesium from contaminated soils. It is a semi-mechanistic model

253    which considers the partitioning of radiocaesium between the clay and

254    humic fractions of soils; the time-dependent fixation of radiocaesium to

255    clay particles; and competition between radiocaesium and potassium ions

256    for plant uptake. The input variables for the model are the physical and

257    chemical characteristics of the contaminated soils, namely: pH, fractional

258    clay content, fractional organic matter content, the radiocaesium activity

259    concentration and the concentrations of exchangeable potassium and

260    ammonium in the soil. The model is schematically presented in Figure 1,

261    which shows the extensive inter-connection between the model's

262    variables, each of which has a specific mechanistic interpretation (Table

263    2).

264    The model was parameterised using data from two comparable

265    experiments in which radiocaesium uptake by grass was measured for a

266    wide range of soil types. The study by Smolders *et al.* (1997) focussed on

267    mineral soils (with relatively low radiocaesium uptake), whereas the study

268    by Sanchez *et al.* (1999) considered organic soils (with relatively high

269    radiocaesium uptake). Employing the definitions given above, the model

270    comprises 6 input variables, 17 model variables, 8 fixed parameters and 7

271    adjustable parameters. The adjustable parameters were estimated by

272    fitting the model to the combined data set using the Marquardt non-linear

273   regression method (Press *et al*., 1989). An additional data set, derived

274   from the work of Nisbet *et al.* (1999), provided an independent test of the

275   model's predictive performance. This data provided sufficient information

276   for the application of the model, although it considers a range of

277   graminaceous cereals rather than grass specifically. Consequently, it

278   might be expected to show a higher degree of variability than the data set

279   used to fit the models (the parameterisation data set).

280

281   Implementation

282   The original model was run using the full range of soil input variables

283   within Absalom *et al*.'s (2001) parameterisation data, to allow the mean

284   values of the model variables to be calculated.

285   As a preliminary screening procedure all the model variables were

286   individually replaced (i.e. with all other variables retaining their original

287   formulation) to identify potential replacement candidates. Any model

288   variable whose replacement did not more than double the RSS with

289   respect to the parameterisation dataset was deemed a replacement

290   candidate. This procedure identified 10 model variables: ph, $M_{CaMg}$, $CEC_h$,

291   $CEC_c$, $\theta_h$, $Kx_s$, $NH_4$, $Kd_h$, $\theta_c$ and $RIP_c$. An exhaustive simplification was then

292   performed, whereby a model formulation was generated for every possible

293   combination of replacement of these model and input variables ($2^{10}=1024$

294   in total).

295   For each reduced model the adjustable parameters were re-estimated

296   using the Marquardt procedure (Press *et al*., 1989) originally employed by

297   Absalom *et al*. (2001). In each case, the parameterisation data were used

298    to calculate RSS, $AIC_c$, BIC, MDL and ICOMP. The independent data

299    derived from Nisbet *et al*. (1999) were used to calculate the prediction

300    sum of squares (PSS), which was used as an indicator of the model's

301    general predictive capability.

302

303    <u>Results</u>

304    The models with the best performance measures for each criterion are

305    summarised in Table 3. Two measures of model complexity are shown:

306    the number of adjustable parameters (p), which is the traditional measure

307    of complexity of statistical models, and the number of model and input

308    variables (M), which is arguably a more relevant measure of complexity

309    for mechanistic models although not normally considered in statistical

310    model selection.

311    The lowest values of RSS and $AIC_c$ occurred for the same model, in which

312    $M_{CaMg}$, $CEC_h$, and pH were replaced. As can be seen in Figure 1, these

313    three variables are directly related, and replacing pH has the effect of also

314    replacing $CEC_h$ and $M_{CaMg}$ with constants. Similarly, if both $CEC_h$ and $M_{CaMg}$

315    are replaced, pH can effectively be considered a constant. This model had

316    a lower RSS than the full model (36.84 c.f. 39.15). In this case the

317    number of adjustable parameters is the same as in the original model (i.e.

318    7), although the number of model and input variables is reduced from 22

319    to 19. This arises because the replaced variables ($M_{CaMg}$, $CEC_h$, and pH) do

320    not utilise any adjustable parameters (the use of adjustable parameters is

321    indicated in Figure 1).

322   The lowest values of BIC, MDL and ICOMP were all associated with a

323   further reduced model in which $Kd_h$ and $RIP_C$ were replaced, in addition to

324   $M_{CaMg}$, $CEC_h$, and pH. This model had a higher RSS than the original model.

325   However, p is reduced to 5 due to the replacement of the model variable

326   $RIP_c$, which more than compensates for the loss of fit in the calculation of

327   BIC, MDL and ICOMP.

328   Both reduced models resulted in lower values of PSS than the full model,

329   with the RSS-$AIC_c$ selected model slightly outperforming the BIC-MDL-

330   ICOMP selected model; although this difference appears trivial.

331   For each of the criteria, there was little difference between the best

332   performing models and those models with second lowest criteria scores.

333   In all cases, the only difference was the inclusion or exclusion of $Kd_h$

334   (depending on whether it was present in the best model). Furthermore,

335   this replacement had a relatively small effect on the criteria scores. For

336   example, $RSS_p$ increased from 36.84 to 37.63, BIC increased from 69.03

337   to 69.38, MDL increased from 23.98 to 24.07 and ICOMP increased from

338   25.73 to 25.97 for the best and second-best models respectively.

339   The models with the third lowest criteria scores all involved the

340   replacement of $CEC_c$. This resulted in more significant increases in the

341   respective criteria scores.

342

343   <u>Discussion of example application</u>

344   The two reduced models selected both had the pH input variable replaced,

345   together with the model variables solely dependent upon it. Although this

346   is a very clear finding across all of the performance criteria it is

347  mechanistically surprising. Many subject specialists would expect pH to be

348  related to plant uptake of radiocaesium. However, these results suggest

349  that the pH input variable is introducing additional variation into the model

350  predictions, which is not accounted for by the relationships that predict

351  the soil solution concentration of Ca and Mg ($M_{CaMg}$) and the cation

352  exchange capacity of the humic fraction ($CEC_h$). This does not imply that

353  pH does not play a role in the uptake of radiocaesium, merely that the pH

354  input variable in this model does not contribute to its predictive capability.

355  Pragmatically, the removal of pH increases the utility of the model, as it

356  reduces the model's input requirements. This is especially important in the

357  case of the Absalom model as it has been applied spatially (Gillett *et al.*

358  (2001)), and pH is a difficult soil parameter to obtain from spatial data

359  sets.

360  The further replacement of $RIP_c$ and $Kd_h$ is recommended by BIC-MDL-

361  ICOMP, notwithstanding the increase in $RSS_P$, as this reduces the number

362  of adjustable parameters. These model variables seek to refine the

363  model's description of Cs adsorption in soils, accounting for the differences

364  between adsorption on mineral and clay surfaces. While these may well be

365  real processes the implication of the BIC-MDL-ICOMP result is that these

366  refinements are over-fitting the model to the parameterisation data.

367  Although, the results of the independent test of the model's predictive

368  performances do not support this conclusion, they do suggest there is

369  very little benefit from the inclusion of these variables.

370

## General Discussion

372 The widely used approach of comparing the predictions of a model to

373 corresponding observed values provides a basis for assessing the

374 performance of the model. However, this is a test without a 'scale' unless

375 there is a comparison *between* different models of the same system.

376 The approach described here provides a method for rapidly generating

377 many alternative model formulations, which may then be compared using

378 various performance measures. Of course, all of the model formulations

379 that are generated are based on the structure of the original model.

380 Clearly, we are not investigating all possible models for a system but a

381 related sub-set. For this reason, we regard the approach as a potentially

382 useful diagnostic, which can be used to inform model formulation, rather

383 than as a method for definitively identifying the best model. For example,

384 in the case of the Absalom model the results suggest specific aspects of

385 the model's formulation that could be re-visited.

386 The importance of expert scientific knowledge when designing mechanistic

387 models remains paramount. However, if models are to be used for

388 predictive purposes it is also important that they have empirical support

389 and are not over-fitted. The proposed approach is potentially valuable in

390 this regard, as useful information can be obtained about the empirical

391 justification of hypotheses contained in a model by comparing the

392 numerous simpler models generated with the full model.

393 The example we have presented here included a formal parameterisation

394 step. The application of AIC, BIC, MDL and ICOMP is dependent on this as

395 they are based on the concept of formally fitted parameters and, in the

396 case of MDL and ICOMP, information about the variances and co-variances

397 of parameter estimates. However, this is not a requirement for the

398 application of the simplification approach. The simple comparisons to

399 observed data could be applied to any model, and the use of a data set

400 truly independent of model development is probably a valuable

401 alternative.

402 A limitation to this approach is that an exhaustive search of all possible

403 combinations of model variable replacements may become

404 computationally prohibitive in situations where there are large numbers of

405 candidate variables for replacement. This would be especially true for

406 models that were computationally intensive in their original form. In such

407 cases, it may be that some form of successive search, analogous to

408 stepwise regression procedures, could be developed.

409 An alternative approach to selecting a best model, which is now commonly

410 used in the case of statistical models, is to average predictions over a

411 class of possible models, weighted in some way by their performance (e.g.

412 Hoeting *et al*. (1999)). This type of method is also applicable to

413 alternative mechanistic model formulations and the proposed approach

414 may provide a means of creating appropriate alternative models.

415

420

## References

422 Absalom, J. P., Young, S. D., Crout, N. M. J., Sanchez, A., Wright, S. M.,
423 Smolders, E., Nisbet, A. F., Gillett, A. G., 2001. Predicting the transfer of
424 radiocaesium to plants using soil characteristics. J. Environ. Radioactiv.,
425 52:31-43.

426

427 Akaike, H., 1973. Information theory and an extension of the maximum
428 likelihood principle. In: Petrov, B. N., Csaki, F. (Editors) Second
429 International Symposium on Information Theory. Akademiai Kiado,
430 Budapest. 267-281.

431

432 Asgharbeygi, N., Langley, P., Bay, S., Arrigo, 2006. Inductive revision of
433 quantitative process models. Ecol. Mod., 194:70-79.

434

435 Atanasova, N., Todorovski, L., Džeroski, S. and Kompare, B., 2006.
436 Constructing a library of domain knowledge for automated modelling of
437 aquatic systems. Ecol. Mod., 194:14-36.

438

439 Bozdogan, H., 2000. Akaike's information criterion and recent
440 developments in information complexity. J. Math. Psych., 44:62-91.

441

442 Brooks, R. J., Semenov, M. A., Jamieson, P. D., 2001. Simplifying Sirius:

443 sensitivity analysis and development of a meta-model for wheat yield

444 prediction. Eur. J. Agron., 14:43-60.

445

446 Burnham, K. P., Anderson, D. R., 2002 (Second edition). Model selection

447 and multimodel inference. Springer, New York.

448

449 Crout, N. M. J., Beresford, N. A., Howard, B. J., Mayes, R. W., Hansen, H.

450 S., 1998. A model of radiostrontium transfer in dairy goats based on

451 calcium metabolism. J. Dairy Sci., 81:92-99

452

453 Gibbons, J. M., Sparkes, D. L., Wilson, P., Ramsden, S.J., 2005. Modelling

454 optimal strategies for decreasing nitrate loss with variation in weather – a

455 farm-level approach. Agr. Syst., 83:113-134.

456

457 Gillett, A. G., Crout, N. M. J., Absalom, J. P., Wright, S. M., Young, S. D.,

458 Howard, B. J., Barnett, C. L., McGrath, S. P., Beresford, N. A., Voigt, G.,

459 2001. Temporal and spatial prediction of radiocaesium transfer to food

460 products. Radiat. Environ. Biophys., 40:227-235.

461

462 Hoeting, J. A., Madigan, D., Raftery, A. E., Volinsky, C. T., 1999. Bayesian

463 model averaging: A tutorial. Stat. Sci., 14:382-401.

464

465    Hurvich, C. M., Tsai, C-L., 1989. Regression and time series model

466    selection in small samples. Biometrika, 76:297-307.

467

468    Jamieson, P. D., Semenov, M. A., Brooking, I. R. Francis, G. S., 1998.

469    Sirius: a mechanistic model of wheat response to environmental variation.

470    Eur. J. Agron., 8:161-179.

471

472    Myung, J., 2000. The importance of complexity in model selection. J.

473    Math. Psych., 44:190-204.

474

475    Myung, J., Pitt, M. A., 2002. When a good fit can be bad. Trends. Cogn.

476    Sci., 6:421-425.

477

478    Nisbet, A. F., Woodman, R. F. M., Haylock, R. G. E., 1999. Recommended

479    soil-to-plant transfer factors for radiocaesium for use in arable systems.

480    NRPB-R304. National Radiological Protection Board, Chilton, Didcot, UK.

481

482    Press, W. H., Flannery, B. P., Teukolsky, S. A., Vetterling, W. T. 1989.

483    Numerical recipes in Pascal. Cambridge University Press, Cambridge, UK.

484    Raftery AE. 1995. Bayesian model selection in social research. Sociological

485    Methodology 25:111-163.

486

487   Rissanen, J., 1987. Stochastic complexity and the MDL principle.

488   Econometric Reviews, 6:85-102.

489

490   Sanchez, A. L., Wright, S. M. Smolders, E., Naylor, C. Stevens, P. A.,

491   Kennedy, V. H., Dodd, B. A., Singleton, D. L., Barnett, C. L., 1999. High

492   plant uptake of radiocaesium from organic soils due to Cs mobility and low

493   soil K content. Environ. Sci. Technol., 33:2752-2757.

494

495   Schwarz, G., 1978. Estimating the dimension of a model. Ann. Statist., 6:

496   461-464.

497

498   Smolders, E., Van Den Brande, K., Merckx, R., 1997. The concentrations

499   of $^{137}$Cs and K in soil solution predict the plant availability of $^{137}$Cs in soils.

500   Environ. Sci. Technol., 31:3432-3438.

501 **Tables**

502 Table 1. Commonly used model selection criteria.

| Criterion | Calculation | | Reference |
|---|---|---|---|
| | GOF term | Complexity term | |
| AIC | -2ln(ML) + | 2p | Akaike (1973) |
| $AIC_c$ | -2ln(ML) + | 2p+ 2p(p+1)/(n-p-1) | Hurvich and Tsai (1989) |
| BIC | -2ln(ML) + | p*ln(n) | Schwarz (1978) |
| MDL | -ln(ML) + | ½ln(\|H\|) | Rissanen (1987) |
| ICOMP | -ln(ML) + | (p/2)ln(tr($\theta$)/p) – ½ ln\|$\theta$\|) | Bozdogan (2000) |

503 Where: ML is the maximised likelihood; p is the number of parameters
504 estimated using data; n is the number of data points used to determine
505 the maximum likelihood; H is the Hessian matrix; tr($\theta$) is the trace of the
506 parameter covariance matrix.

507

508

509

510

511

512

513

514    Table 2. Mechanistic descriptions of variables in the Absalom model.

| Model variable | Mechanistic interpretation | Units/scale |
|---|---|---|
| % clay | Fraction of clay matter in soil | % |
| % C | Fraction of organic matter in soil | % |
| $K^+$ | Exchangeable potassium in soil | Meq $100g^{-1}$ |
| pH | Soil pH | 0-14 |
| $NH_4$ | Ammonium concentration in soil | Mol $dm^{-3}$ |
| $\theta_c$ | Gravimetric clay content | g $g^{-1}$ |
| $\theta_c$ | Gravimetric clay content | g $g^{-1}$ |
| $RIP_c$ | Radiocaesium interception potential | mmol $kg^{-1}$ |
| $Kx_{soil}$ | Exchangeable potassium in soil | $Cmol_c$ $kg^{-1}$ |
| $CEC_h$ | Cation exchange capacity on the humic soil fraction | $Cmol_c$ $kg^{-1}$ |
| M_camg | Concentration of Calcium and Magnesium ions in the soil solution | Mol $dm^{-3}$ |
| $CEC_c$ | Cation exchange capacity on the clay soil fraction | $Cmol_c$ $kg^{-1}$ |
| $Kx_h$ | Exchangeable potassium on the humic soil fraction | $Cmol_c$ $kg^{-1}$ |
| $Kd_h$ | Radiocaesium distribution coefficient for the humic soil fraction | mol $kg^{-1}$ |
| $Kd_c$ | Radiocaesium distribution coefficient for the clay soil fraction | mol $kg^{-1}$ |
| mk | Concentration of $K^+$ in the soil solution | mol $dm^{-3}$ |
| Kdr | Proportion of labile $Cs^+$ adsorbed on the clay fraction | 0-1 |
| Kdl | Total labile radiocaesium | mol $kg^{-1}$ |
| CF | Concentration factor | $dm^3$ $kg^{-1}$ |
| D factor | Dynamic factor which describes the change in labile $Cs^+$ with time. | 0-1 |
| $Cs_{sol}$ | Radiocaesium activity concentration in soil solution | Bq $dm^{-3}$ |
| $Cs_p$ | Radiocaesium activity concentration in plants | Bq $kg^{-1}$ |
| $Cs_{soil}$ | Total radiocaesium activity concentration in soil | Bq $kg^{-1}$ |

Table 3.  Summary of the original model and the best performing reduced models selected by RSS, $AIC_c$, BIC, MDL and ICOMP.

| Selection criterion | Model variable | | | | | | | | | | p | M | RSS | PSS |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $M_{camg}$ | $CEC_h$ | $NH_4$ | $CEC_c$ | pH | $\theta_h$ | $Kx_s$ | $\theta_c$ | $Kd_h$ | $RIP_c$ | | | | |
| None (full model) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 7 | 22 | 39.15 | 20.69 |
| RSS, $AIC_c$ | × | × | ✓ | ✓ | × | ✓ | ✓ | ✓ | ✓ | ✓ | 7 | 19 | 36.84 | 16.59 |
| BIC, MDL, ICOMP | × | × | ✓ | ✓ | × | ✓ | ✓ | ✓ | × | × | 5 | 17 | 43.69 | 16.68 |

✓ indicates that the variable remains in the model in its original form and × denotes that the variable is replaced by a constant. RSS is the residual sum of squares for the parameterisation dataset; PSS is the prediction sum of squares for the independent dataset; *p* indicates the number of adjustable parameters present in the model; M indicates the number of model and input variables in the model.