

# EM and component-wise boosting for Hidden Markov Models: a machine-learning approach to capture-recapture

Robert W. Rankin<sup>a,\*</sup>

<sup>a</sup>*Cetacean Research Unit, School of Veterinary and Life Sciences, Murdoch University, Australia*

---

## 1 Abstract

2 This study presents a new boosting method for capture-recapture models, rooted in predictive-  
3 performance and machine-learning. The regularization algorithm combines Expectation-Maximization and  
4 boosting to yield a type of multimodel inference, including automatic variable selection and control of model  
5 complexity. By analyzing simulations and a real dataset, this study shows the qualitatively similar estimates  
6 between AICc model-averaging and boosted capture-recapture for the CJS model. I discuss a number of  
7 benefits of boosting for capture-recapture, including: i) ability to fit non-linear patterns (regression-trees,  
8 splines); ii) sparser, simpler models that are less prone to over-fitting, singularities or boundary-value esti-  
9 mates than conventional methods; iii) an inference paradigm that is rooted in predictive-performance and  
10 free of p-values or 95% confidence intervals; and v) estimates that are slightly biased, but are more stable over  
11 multiple realizations of the data. Finally, I discuss some philosophical considerations to help practitioners  
12 motivate the use of either prediction-optimal methods (AIC, boosting) or model-consistent methods. The  
13 boosted capture-recapture framework is highly extensible and could provide a rich, unified framework for  
14 addressing many topics in capture-recapture, such as spatial capture-recapture, individual heterogeneity, and  
15 non-linear effects.

16  
17 *Keywords: capture-recapture, boosting, machine-learning, model-selection, marked animals, high-dimensional*  
18 *data*

## 19 1. Introduction

20 In this study, I introduce boosting for Hidden-Markov Models (HMM) with a particular focus on capture-  
21 recapture models. It is targeted at capture-recapture practitioners who desire model parsimony under low-  
22 sample sizes and high-dimensional settings. Capture-recapture systems are perennially in an situation of  
23 high model-uncertainty (Johnson & Omland, 2004) and would benefit from an inference-paradigm that is  
24 flexible, extensible and rooted in good *predictive performance*. Some questions are the following. Can we find  
25 a simple model out of the hundreds or millions of plausible “fixed-effects” models? Can we correctly identify

---

\*Corresponding author. E-mail: robertw.rankin@gmail.com

26 a sparse set of highly influential covariates in high-dimensional situations? Can the method accommodate  
27 non-linear relationships and interactions (e.g., regression trees, kernels and splines) without over-fitting? Can  
28 the method avoid the scourge of singularities and boundary-value estimates that trouble MLE-based models  
29 and their model-averaged derivatives? How does the method compare to other popular multimodel inference  
30 techniques, such as AICc model-averaging?

31 A motivating model will be the Cormack-Jolly-Seber (CJS) capture-recapture model, with a focus on  
32 which covariates influence the survival of an open population of marked animals under imperfect detection.  
33 While there are many regularization and variable selection techniques in univariate regression models, the  
34 problem becomes combinatorially difficult for HMMs such as capture-recapture models: we must consider  
35 multiple plausible specifications for both the transition process (survival), as well as the emission process  
36 (capture probability).

37 The issues of model selection and multimodel inference are front-and-centre in most capture-recapture  
38 studies. For example, the popular Program MARK (White & Burnham, 1999) is strongly allied to the  
39 model-averaging ideas of Burnham, Anderson, Buckland and others (Buckland et al., 1997; Anderson et al.,  
40 2000; Burnham, 2004; Burnham & Anderson, 2014). By default, the program offers AICc-weighted averages  
41 (Akaike, 1974) of survival and capture probability. The widespread use of model-averaging in the capture-  
42 recapture field reflects an early appreciation by researchers for the *model uncertainty* inherent to capture-  
43 recapture: every analysis has dozens or thousands of plausible fixed-effect models, including, at a minimum,  
44 time-varying vs time-invariant processes. However, such *post-hoc* model-selection and/or averaging become  
45 computationally unfeasible with just a few extra covariates, due to the combinatorial explosion in the number  
46 of plausible models. Secondly, even if one could realistically compute every model, the AIC/AICc tends to  
47 favour more complex models (Shao, 1997; Hooten & Hobbs, 2015), which, in a capture-recapture context,  
48 can have singularities or boundary-value estimates (like 100% survival or 100% capture probability; Rankin  
49 et al., 2016; Hunt et al., 2016). This latter problem is rarely appreciated, but has motivated the development  
50 of Bayesian models to encourage parsimony under sparse data (Schofield et al., 2009; Schofield & Barker,  
51 2011; Rankin et al., 2014, 2016)

52 Clearly, methods are needed to address the dual challenge of variable selection and low-sample sizes. Also,  
53 we should favour flexible techniques that can accommodate different functional forms (such as regression trees,  
54 splines, random effects) and find covariate-interactions, without over-fitting or producing boundary-value  
55 estimates.

56 Hand & Vinciotti (2003) and Burnham (2004) hinted at a possible contender to the model-averaging ap-  
57 proach when they suggested a parallel between multimodel inference and boosting: whereas model-averaging  
58 weights dozens or hundreds of fixed-effect models, boosting sequentially combines hundreds or thousands of  
59 simple *weak learners* to yield a strong statistical model in aggregate. Most ecologists are familiar with boost-  
60 ing for univariate regression and classifications tasks (Elith et al., 2008; Kneib et al., 2009; Opper et al., 2009;  
61 Hothorn et al., 2010; Tyne et al., 2015), but the recently developed *component-wise boosting* and *gamboostLSS*

62 algorithms (Bühlmann & Yu, 2003; Schmid & Hothorn, 2008b; Schmid et al., 2010; Mayr et al., 2012) opened  
63 the way for complex hierarchical distributions with many components (Hothorn et al., 2010; Hutchinson  
64 et al., 2011; Schmid et al., 2013; Hofner et al., 2014). Under this boosting framework, each boosting iteration  
65 alternates between fitting the capture probability parameter (conditional on survival), and then fitting the  
66 survival component (conditional on the capture probabilities). Plus, boosting offers a wide variety of possible  
67 weak learners, from ordinary least squares to splines and CART-like trees (Hothorn et al., 2006; Bühlmann &  
68 Hothorn, 2007). This gives boosting much appeal over other sparsity-inducing variable selection paradigms,  
69 such as the Lasso (Tibshirani, 2011; Efron et al., 2004), Elastic-Net, Support Vector Machines, Hierarchical  
70 Bayesian shrinkage-estimators (Rankin et al., 2016). In this way, component-wise boosting offers a unified  
71 framework to address high-dimensional variable selection, interaction-detection, and non-linear relationships,  
72 while encouraging model parsimony through a prediction-optimized control on model complexity.

73 The contribution of this study is to develop a boosting method suitable for the Cormack-Jolly-Seber  
74 capture-recapture model (hereafter, CJSboost) and whose framework can be used for a wider class of capture-  
75 recapture models. The particular challenge of boosting a HMM is the serially dependent nature of observa-  
76 tions. Hitherto, boosting methods required independent data points in order to perform gradient descent,  
77 e.g., by descending the point-wise *negative gradient* of a loss-function. The CJSboost approach is to garner  
78 such conditional independence by imputing expectations of *latent states*  $\mathbf{z}$  (here, alive or dead). In CJSboost,  
79 we alternate between boosting the parameters (conditional on latent states) and imputing expectations of the  
80 latent states (conditional on the parameters). I provide two different techniques to impute such expectations:  
81 i) Expectation-Maximization (called CJSboost-EM), and ii) Monte-Carlo approximation of the marginal dis-  
82 tribution of latent states (CJSboost-MC). As I will show, both algorithms lead to approximately the same  
83 estimates. Furthermore, the estimates are qualitatively very similar to the model-averaged estimates by AICc  
84 weighting. The AIC is also motivated by optimal (asymptotic) predictive performance.

85 This article will demonstrate CJSboost via simulations and an analysis of an European Dipper dataset  
86 from Lebreton et al. (1992), with particular emphasis on comparing estimates from linear and non-linear mod-  
87 els (e.g., CART-like trees), and comparisons to Maximum Likelihood estimation and AICc model-averaging  
88 (Burnham, 2004) using Program MARK (White & Burnham, 1999). Simulations will also challenge CJS-  
89 boost to perform a model-selection task that is nearly impossible for conventional methods: finding a sparse  
90 set of influential covariates among  $21 \times 21$  different covariates.

91 There are two potential audiences for this paper. First, HMM practitioners will be interested in a  
92 general approach to boosting and HMMs, which opens new possibilities for incorporating regularization, semi-  
93 parametric learners and interaction detection to a vast catalogue of applications. For the second audience of  
94 mark-recapture practitioners, I offer a fresh view of mark-recapture from a *prediction* or *learning* perspective.  
95 For example, we can observe the degree to which regularization and bootstrap-validation suggest simpler  
96 models than those implied by AICc model-averaging. Boosting also offers capture-recapture an alternative  
97 means of inference that is principled and free of p-values and 95% Confidence Intervals (Anderson et al.,

98 2000; Hoekstra et al., 2014; Morey et al., 2016). Furthermore, this new capture-recapture paradigm can  
99 easily accommodate a range of hot-topics in capture-recapture, such as individual-heterogeneity and spatial  
100 capture-recapture, by leveraging the wide variety of base-learners available in the `mboost` family of R packages  
101 (Bühlmann & Hothorn, 2007; Hothorn et al., 2006; Mayr et al., 2012; Hofner et al., 2012).

## 102 2. Methods

### 103 2.1. Organization

104 The manuscript begins by introducing some basic ideas of statistical learning theory (Section 2.2) and the  
105 Cormack-Jolly-Seber model. Section 2.3 describes two boosting algorithms, CJSboost-EM and CJSboost-  
106 MC, for capture-recapture models. Section 2.4 discusses some important practical considerations about  
107 regularization and base-learners. Section 2.5 describes a simulation to compare the estimates from CJSboost-  
108 EM and CJSboost-MC, as well as AICc model-averaging and MLEs (results in 3.1). Section 2.6 describes  
109 a reanalysis of of dipper dataset using CJSboost-EM and AICc model-averaging (results in 3.2). Section  
110 2.7 uses simulations to assess the performance of CJSboost-EM under a high-dimensional model-selection  
111 problem (results in 3.3). The manuscript finishes with a discussion about how to interpret the results from  
112 CJSboost and poses some new questions (Section 4). A summary is provided in Section 5. For R code and  
113 a tutorial, see the online content at <http://github.com/faraway1nspace/HMMboost/>.

### 114 2.2. Background

#### 115 2.2.1. The Prediction perspective

From a prediction perspective, our goal is to estimate a prediction function  $G$  that maps covariate in-  
formation  $\mathbb{X}$  to our response variable (i.e.,  $G: \mathbb{X} \rightarrow \mathbb{Y}$ ). Our data  $\{y_j, \mathbf{x}_j\}_{j=1}^n$  arises from some unknown  
probability distribution  $P$ . Our optimal prediction function is that which minimizes the *generalization error*:

$$\mathcal{L}(y, G(\mathbf{x})) = \int \ell(y, G(\mathbf{x})) dP(y, \mathbf{x}) \quad (1)$$

116 where  $\ell$  is a *loss* function (it scores how badly we are predicting  $y$  from  $\mathbf{x}$ ) and  $\mathcal{L}$  is the *expected loss*, a.k.a,  
117 the *risk* (our loss integrated over the entire data distribution). Our goal is to minimize the loss on new,  
118 unseen data drawn from the unknowable data distribution  $P$  (Bühlmann & Yu, 2003; Meir & Rätsch, 2003;  
119 Murphy, 2012a). It should be noted that for many disciplines, making good predictions is the primary goal  
120 (e.g., financial forecasting). In mark-recapture, we usually wish to make inference about covariates  $\mathbb{X}$  and  
121 their functional relationship ( $G$ ) to the response variable, such as estimating survival from capture histories,  
122 rather than making predictions *per se*. In such cases, the generalization criteria (1) serves as a principled  
123 means of “model parsimony”: our model is as complex as is justified to both explain the observed data  
124 and make good predictions on new data. This is very different from Maximum-Likelihood Estimation (as in  
125 Program MARK) whose estimate  $\hat{G}$  is that which maximizes the likelihood of having seen the observed data  
126  $\mathbf{y}$ . It is, however, similar to AIC selection, which is implicitly motivated by minimizing expected loss (Vrieze,  
127 2012), i.e., optimal predictive performance.

One cannot measure the generalization error (1); instead, we must proceed by minimizing the *empirical risk* measured on our observed data:

$$L(\mathbf{y}, G(\mathbf{X})) = \sum_{j=1}^n \ell(y_j, G(\mathbf{X}_j)) \quad (2)$$

128 Minimizing  $L(\mathbf{y}, G(\mathbf{X}))$  until convergence is easy but will obviously *over-fit* a sample and make bad predic-  
129 tions. However, it can be shown that if we constrain the complexity of our function space (Bühlmann & Yu,  
130 2003; Meir & Rätsch, 2003; Mukherjee et al., 2003) we can pursue a strategy of “regularized risk minimization”  
131 and bound the generalization error. In learning algorithms, this entails at least one regularization parameter  
132 that smooths or constrains the complexity of  $G$ . In other words, we do not seek the estimator that best fits  
133 the data. In boosting, the principal means of regularization is via shrinkage (taking only small steps along the  
134 risk gradient) and early-stopping (*not* running the algorithm until the risk convergences). These correspond  
135 to hyperparameters  $\nu$  and  $m$ , respectively, the shrinkage weight and the number of boosting iterations. For  
136 a small  $m$ , the model is strongly constrained and very conservative; as  $m$  gets big, the model becomes more  
137 complex. Likewise,  $\nu \ll 1$  ensures that the influence of any one boosting step is tiny. Practically, one fixes  $\nu$   
138 and finds an optimal  $m$  via cross-validation. Figure 1 shows an example of bootstrap-validation to find an  
139 optimal stopping criteria  $m_{CV}$ , used in the dipper CJS analysis (Section 3.2).

#### 140 2.2.2. Motivation for regularization

141 The unregularized boosted model with prediction function  $G^{(m \rightarrow \infty)}$  results in a *fully-saturated model*,  
142 which (depending on the prediction function) is equivalent to the Maximum Likelihood solution (Mayr et al.,  
143 2012). At finite sample sizes and a large candidate-set of covariates, MLEs do not result in good predictions:  
144 they may be unbiased, but they will be wildly sensitive to noise in the data, especially for capture-recapture.  
145 For a regularized model  $G^{(m \ll \infty)}$ , learning algorithms should preferentially select influential covariates and  
146 shrink the coefficients of less-important covariates close to zero. This shrinkage induces a bias (Bühlmann &  
147 Yu, 2003; Bühlmann & Hothorn, 2007), but the predictions are more robust to noisy data (i.e. low-variance;  
148 Murphy, 2012a). In this light, we see the practical similarity between regularization and the more popular  
149 model parsimony strategies in capture-recapture, such as model-selection, model-averaging, and subjective  
150 Bayesian models. Hooten & Hobbs (2015) implore ecologists to unify these techniques under a Bayesian  
151 perspective; for example, the AIC, Lasso/L2Boosting, Ridge-regression can be reformulated in such a way  
152 that they differ according to the priors on the  $\ell_0$ ,  $\ell_1$  and  $\ell_2$ -norm of regression-coefficients, respectively. Even  
153 a simple Bayesian prior can be understood as a type of regularization by shrinking estimates away from  
154 their MLE values and towards the conservative expectations of a prior (a.k.a “natural shrinkage”; Hooten &  
155 Hobbs, 2015).

156 Today, most mark-recapture practitioners are implicitly using a prediction criteria for inference. For  
157 example, the AIC is popular in mark-recapture studies (Johnson & Omland, 2004), thanks in large part  
158 to the Frequentist and Information-Theoretic leanings of Program MARK (White & Burnham, 1999). The  
159 AIC is asymptotically prediction-optimal, whose maximum risk is minimal among all potential models, and

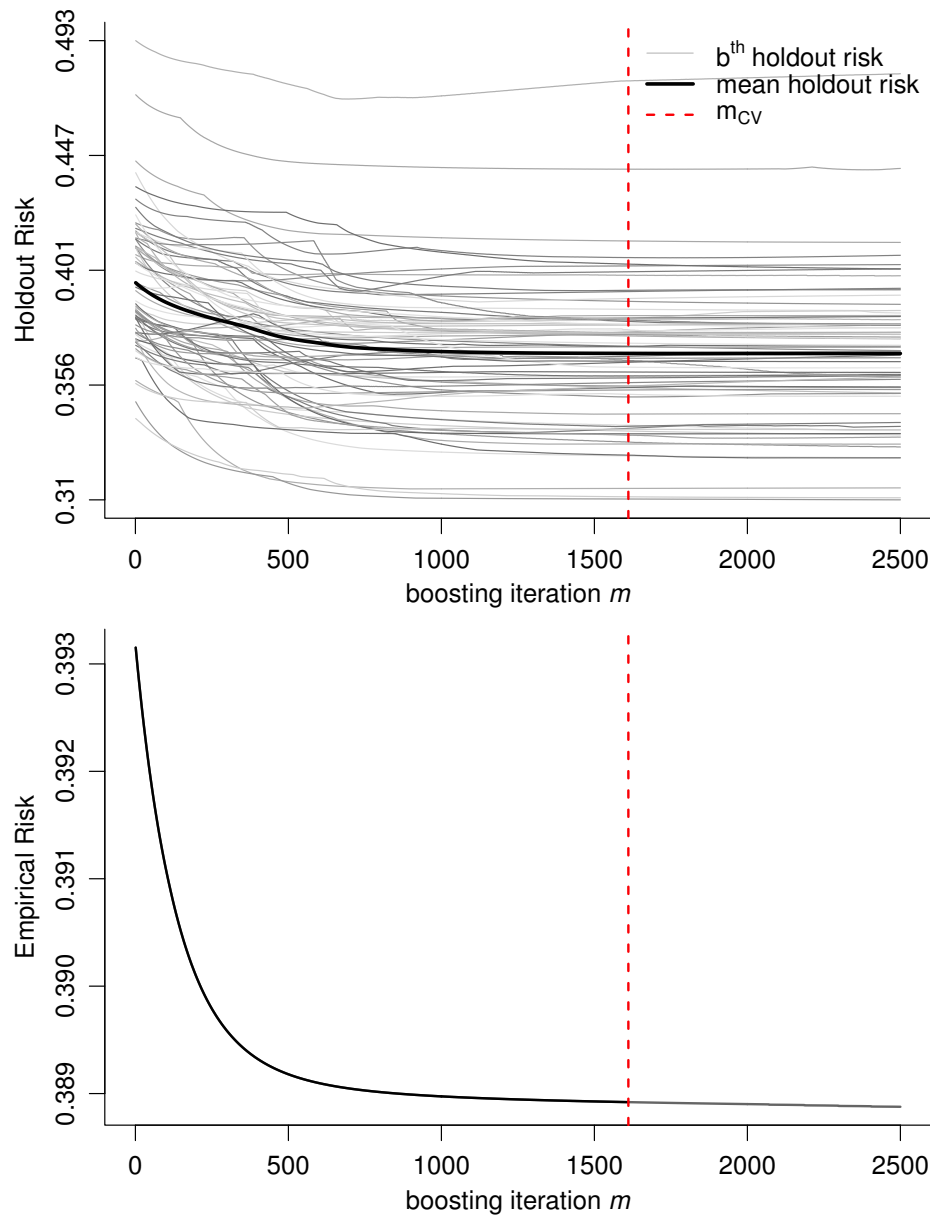


Figure 1: Monitoring the risk minimization (negative CJS log-Likelihood) for the Dipper analysis, using CJSboost-EM with CART-like trees as base-learners. Each boosting iteration  $m$  takes a step towards minimizing the empirical risk and selects a new shrunken base-learner to add to the ensembles. *Top*: Estimating the optimal stopping criteria at  $m_{\text{CV}}$  (red dashed line) via bootstrap-validation. Each grey line represents the holdout-risk predicted on a subset of the capture-histories, from a model trained on a bootstrap-sample of capture-histories.  $m_{\text{CV}}$  minimizes the mean holdout-risk over all bootstrap runs, an estimate of the *expected loss*. *Bottom*: The empirical risk of the final statistical model using the full dataset; stopping early at  $m_{\text{CV}}$ , well before convergence.

160 has connections with leave-one-out-cross-validation (LOOCV; Stone, 1977; Shao, 1993, 1997; Vrieze, 2012).  
 161 However, statisticians consider the AIC to be a bit too permissive, especially if the “true model” is sparse  
 162 (Shao, 1993; Burnham, 2004; Hooten & Hobbs, 2015). For practical mark-recapture analysis, the AIC/AICc  
 163 can favour overly-complex models which can suffer singularities or boundary-value estimates (like 100%  
 164 survival or 100% capture probability; Rankin et al., 2016; Hunt et al., 2016). Boosting is also prediction-

165 optimal (Bühlmann & Yu, 2003), but skirts the issues of singularities and boundary-value estimates by fitting  
166 very simple models, called *base-learners* in a step-wise manner. At finite sample sizes, boosting should lead  
167 to slightly sparser models than the AIC/AICc.

168 In an extreme case of sparsity, when being prediction-optimal is not the chief concern, and one wishes to  
169 instead uncover a “true model” with just a few important covariates, boosting has another desirable prop-  
170 erty. Regularized risk-minimizers (in a univariate setting) can be made model-selection consistent by hard-  
171 thresholding unimportant covariates to zero weight (Bach, 2008; Meinshausen & Bühlmann, 2010; Murphy,  
172 2012c). These sparse solutions may be more interesting for capture-recapture practitioners when inference  
173 about covariates or estimating survival is the chief concern.

### 174 2.2.3. Introduction to boosting

175 Boosting is an iterative method for obtaining a statistical model via gradient descent (Breiman, 1998;  
176 Friedman et al., 2000; Friedman, 2001; Breiman, 1999; Schmid et al., 2010; Robinzonov, 2013). The key  
177 insight is that one can build a strong predictor  $F = G(X)$  by the step-wise addition of many weak *base-*  
178 *learners*,  $b(y, x) \Rightarrow g$ ,  $g(x): x \rightarrow y$  (Schapire, 1990; Kearns & Valiant, 1994). Remarkably, a base-learner  
179 need only have a predictive performance of slightly better than random chance for the entire ensemble to be  
180 strong. The ensemble results in a smooth additive model of adaptive complexity:

$$F^{(m)} = G(\mathbf{X}) = \sum_{m=1}^{m_{\text{stop}}} \nu \cdot g_k^{(m)}(\mathbf{X}_k) \quad (3)$$

181 where each  $g_k$  is a base-learner’s prediction function, shrunk by  $\nu$ . The ensemble is constructed as follows:  
182 i) initialize the prediction vector  $F^{(m=0)}$  at some uniform estimate (like the MLE of an intercept model); ii) fit  
183 base-learners  $b$  to  $\hat{\mathbf{u}}$ , the estimated negative-gradient of the loss function (the residual variation unexplained  
184 by  $F^{(m-1)}$ ),  $b(\hat{\mathbf{u}}, \mathbf{x}) \Rightarrow g$ ; iii) shrink each base-learners’ prediction  $g(\mathbf{x}) = \hat{f}^{(m)}$  by a small fraction  $\nu$ ; iv)  
185 update the overall prediction  $F^{(m)} = F^{(m-1)} + \nu \hat{f}^{(m)}$ ; v) repeat for  $m_{\text{stop}}$  iterations.  $m_{\text{stop}}$  is the key  
186 parameter that governs model complexity (Bühlmann & Yu, 2003; Schmid & Hothorn, 2008a) and must be  
187 tuned by cross-validation or bootstrap-validation. Variable selection can be directly embedded within each  
188 boosting iteration by choosing only one best-fitting base-learner per  $m$  iteration, discriminating among a  
189 large candidate set of base-learners  $\{b(\mathbf{u}, \mathbf{x}_1), b(\mathbf{u}, \mathbf{x}_2), \dots, b(\mathbf{u}, \mathbf{x}_k)\}$ , and where each candidate only includes  
190 a small subset of the covariates  $\mathbf{X}$ . For linear base-learners, this boosting algorithm is generally considered  
191 equivalent to  $\ell_1$  regularization (Efron et al., 2004; Bühlmann & Hothorn, 2007), a.k.a the Lasso.

192 Base-Learners may be simple Least-Squares estimators,  $b_{\text{OLS}}$ , in which case an unregularized boosted  
193 model will estimate regression coefficients that are practically identical to a frequentist GLM. However,  
194 Bühlmann & Yu (2003) showed that for L2Boosting, good overall predictive performance depends on the fact  
195 that base-learners are very weak. Therefore, practitioners commonly use highly-constrained base-learners  
196 such as Penalized Least Squares  $b_{\text{PLS}}$ , or recursive-partitioning trees  $b_{\text{trees}}$  (a.k.a CART), or low-rank splines  
197  $b_{\text{spline}}$ . Despite their weakness, Bühlmann & Yu note that for a fixed constraint (such as low degrees-of-  
198 freedom in  $b_{\text{spline}}$  or low tree-depth in  $b_{\text{trees}}$ ), the overall boosted ensemble will typically have a much greater



199 complexity than its constituent base-learners and that this complexity is adaptive.

200 There are many flavours of boosting. CJSboost hails primarily from the component-wise boosting and  
 201 gamboostLSS frameworks (Bühlmann & Yu, 2003; Schmid & Hothorn, 2008b; Schmid et al., 2010; Mayr  
 202 et al., 2012). Here, the prediction vector is now a set  $\mathcal{F} = (F_1, F_2, \dots, F_k)$  of  $k$  components, each representing  
 203 one of the parameters in the likelihood function (e.g.,  $\phi$  and  $p$ ). Each parameter has its own ensemble of base-  
 204 learners. The loss function is the negative log-likelihood of the data-generating model  $\ell_i = -\log p(\mathbf{y}_i | \phi_i, p_i) =$   
 205  $-\log p\left(\mathbf{y}_i \mid \frac{1}{1+e^{-F_\phi}}, \frac{1}{1+e^{-F_p}}\right)$  (see the CJS likelihood 4). Each components' gradient can be estimated from  
 206 the negative partial-derivatives of the loss function with respect to  $F_k$ , i.e.,  $\hat{u}_{k,i} = -\frac{\partial \ell_i}{\partial F_k}$ , conditional on the  
 207 values of the other prediction vectors  $F_{-k}$ . Each  $k$  parameter is updated once per boosting iteration.

#### 208 2.2.4. The Cormack-Jolly-Seber model and Hidden Markov Models

209 The above component-wise boosting framework is not suitable for serially dependent observations in an  
 210 HMM time-series: consider that the negative gradient in traditional boosting must be estimated point-wise  
 211 for each independent data pair  $(y_i, X_i)$ . Instead, the CJS likelihood is evaluated on individuals' entire *capture*  
 212 *histories*  $\mathbf{y}_i = (y_1, y_2, \dots, y_T)^\top$  over  $T$  capture periods:

$$p(\mathbf{y}_i | \phi, p, t_i^0) = \left( \prod_{t>t_i^0}^{t_i^*} \phi_{i,t} p_{i,t}^{y_{i,t}} (1 - p_{i,t})^{1-y_{i,t}} \right) \chi_i^{(t_i^*+1)} \quad (4)$$

213 Where  $i$  indexes the  $n$  uniquely identified individuals constituting our dataset;  $t = 1:T$  indexes the  $T$  equally  
 214 spaced capture periods (time);  $y_{i,t} \in [0, 1]$  scores whether individual  $i$  was observed in capture period  $t$ ;  $\phi_{i,t}$  is  
 215 the probability of surviving from capture period  $t-1$  to  $t$  (note the one-time-step difference from the definition  
 216 of  $\phi_t$  used in Program MARK);  $p_{i,t}$  is the capture probability of individual  $i$  in capture period  $t$  (a.k.a, our  
 217 observation error, or the “emission process” in HMM parlance);  $t_i^0$  is the first capture period in which  
 218 individual  $i$  was first observed;  $t_i^*$  is the last period when individual  $i$  was observed. Finally,  $\chi_i^{(t_i^*+1)}$  is the  
 219 probability of never being seen again after  $t_i^*$  until the end of the study,  $\chi_i^{(t)} = (1 - \phi_{i,t}) + (1 - p_{i,t})\phi_{i,t}\chi_i^{(t+1)}$ ,  
 220 and whose recursive calculation exemplifies the serially dependent nature of the model.  $\mathbf{Y}^{n \times T}$  is the full  
 221 matrix of our capture-histories.

222 Mark-recapture practitioners will be interested to note: i) the model conditions on first-capture  $\{t_i^0\}_{i=1}^n$ ;  
 223 ii) the model can potentially allow for individual heterogeneity in capture probabilities  $p_{t,i}$  (which otherwise  
 224 results in a negative-bias in population abundance estimates; Carothers, 1973; Burnham & Overton, 1978;  
 225 Rankin et al., 2016); and iii) certain parameters cannot be separated in Maximum-Likelihood Estimation,  
 226 such as  $p_T$  and  $\phi_T$ , but this is less of an issue under constrained base-learners and regularization.

227 In order to boost the CJS model, we need independence of data pairs  $(y_{i,t}, X_{i,t})$ . If we reformulate the  
 228 capture-recapture system as a HMM, we can garner conditional independence via the concept of latent states  
 229  $z_{i,t} \in \{0, 1\}$  to represent {dead, alive}. When  $z_{i,t} = 1$ , then individual  $i$  is alive and available for capture at  
 230 time  $t$ , and the probability of a capture is simply  $p(y_{i,t}=1 | z_{i,t}=1) = p_{i,t}$ . However, if  $z_{i,t} = 0$  then individual  
 231  $i$  is dead and unavailable for capture at time  $t$ ; therefore the probability of a capture is zero.



232 Obviously, one never knows with certainty the latent states of a trailing sequence of zeros  $\mathbf{y}_{t:T} = (0, \dots, 0)^\top$ ,  
 233 but we can utilize well-developed HMM tools to estimate the state-sequence  $\mathbf{z}$  in various ways. In particular,  
 234 “CJSboost-EM” 2.3.1 utilizes the marginal *expectations* of  $(z_t, z_{t-1})$  in an Expectation-Maximization step.  
 235 “CJSboost-MC” 2.3.2 utilizes Monte-Carlo integration by drawing random values of  $\mathbf{z}$  from the posterior  
 236  $\pi(\mathbf{z}|\mathbf{y}, \phi, p)$ . We can interweave these two methods within a boosting algorithm: both will allow us to  
 237 estimate point-wise negative gradients for all *complete-data* points  $(\{y_{i,t}, z_{i,t}, z_{i,t-1}\}, X_{i,t})$  and proceed with  
 238 the gradient descent algorithm.

### 239 2.3. CJSboost

240 I will now formally describe the CJSboost variants “CJSboost-EM” and “CJSboost-MC”. In practise, I  
 241 will show that they lead to approximately the same estimates, but have different computation disadvantages  
 242 under different scenarios. When the number of discrete states in the HMM process is low (2-3), then the  
 243 deterministic EM algorithm is significantly faster and less prone to approximation error. For example,  
 244 in our CJS example, we just have two latent states  $\{0, 1\} := \{\text{dead}, \text{alive}\}$  with three legal transitions  
 245  $\{1 \rightarrow 1, 1 \rightarrow 0, 0 \rightarrow 0, 0 \rightarrow 1\}$ . However, as the number of discrete states increases, the memory management of  
 246 all the possible transitions becomes combinatorially expensive. In such scenarios, it is computationally easier  
 247 to sample  $z$  from its posterior.

#### 248 2.3.1. CJSboost by Expectation-Maximization

For a CJS model using CJSboost-EM, our target risk is the CJS negative log-likelihood. However, we use the principle of Expectation-Minimization to derive a slightly different loss function and subsequent negative gradients. Our loss is derived from the negative Complete-Data Log-Likelihood (CDL) which assumes we have estimates of the latent state  $z_{i,t}, z_{i,t-1}$ .

$$\begin{aligned}
 -\text{CDL}(y_{i,t}, z_{i,t}, z_{i,t-1} | F_{i,t,\phi}, F_{i,t,p}) = & -\mathbb{1}[z_{i,t-1}=1, z_{i,t}=1] \left( \log \left( \frac{1}{1+e^{-F_{i,t,\phi}}} \right) + y_{i,t} \log \left( \frac{1}{1+e^{-F_{i,t,p}}} \right) \right) \\
 & + (1-y_{i,t}) \log \left( \frac{1}{1+e^{F_{i,t,p}}} \right) \\
 & - \mathbb{1}[z_{i,t-1}=1, z_{i,t}=0] \log \left( \frac{1}{1+e^{F_{i,t,\phi}}} \right) \\
 & - \mathbb{1}[z_{i,t-1}=0, z_{i,t}=0]
 \end{aligned} \tag{5}$$

where  $y$  and  $z$  are defined as above in (4) and  $F_{i,t,p}$  and  $F_{i,t,\phi}$  are the prediction vectors for the capture probability component and the survival component, respectively, on the logit scale. In accordance with the principle of EM, we derive a “Q-function” to serve as our new loss, replacing the values of  $(z_{i,t-1}, z_{i,t})$  with their *two-slice marginal* expectations:  $w_t(q, r) := p(z_{t-1} = q, z_t = r | \mathbf{y}, \mathcal{F})$ . Conditional on the prediction vectors  $\mathcal{F}$  and the capture history  $\mathbf{y}$ , the values of the two-slice marginals  $\{w(1, 1), w(1, 0), w(0, 0)\}$  can be easily computed using a standard “forwards-backwards” HMM algorithm (Rabiner, 1989; Murphy, 2012b), as detailed in Appendix A. We can also treat each  $i \times t$  observation as being conditionally independent,

resulting in the new index  $j := (i, t)$ . The Q-function is:

$$\begin{aligned} \ell(y_j, \{F_{j,\phi}, F_{j,p}\}) &= -w_j(1, 1) \left( \log \left( \frac{1}{1+e^{-F_{j,\phi}}} \right) + y_j \log \left( \frac{1}{1+e^{-F_{j,p}}} \right) + (1-y_j) \log \left( \frac{1}{1+e^{F_{j,p}}} \right) \right) \\ &\quad - w_j(1, 0) \log \left( \frac{1}{1+e^{F_{j,\phi}}} \right) \\ &\quad - w_j(0, 0) \end{aligned} \quad (6)$$

249 According to the theory of EM, by minimizing the Q-function, we also minimize the target empirical  
250 risk: the negative CJS log-likelihood (4). The advantage of working with the Q-function is that it is easy to  
251 calculate the negative gradients (7) and proceed with the gradient descent.

252 I now describe the CJSboost-EM algorithm.

- 253 1. Set the regularization parameters:  $m_{\text{stop}} \approx 10^2 - 10^3$ ;  $\nu_\phi, \nu_p \approx 10^{-3} - 10^{-1}$ ;
- 254 2. Initialize:  $m = 1$ ;  $\mathcal{F}^{(0)} = \{F_\phi^{(0)} = \hat{\phi}^{\text{MLE}}(\cdot), F_p^{(0)} = \hat{p}^{\text{MLE}}(\cdot)\}$  (i.e., initialize the prediction vectors at the  
255 MLEs of a simple intercept model).
- 256 3. Estimate the two-slice marginal probabilities  $\{w_j(1, 1), w_j(1, 0), w_j(0, 0)\}_{j=1}^J$  for all individuals and  
257 capture-periods, using the forwards-backwards algorithm (see Appendix A.3).
4. Estimate the negative gradients:

$$\begin{aligned} \hat{u}_{j,\phi}^{(m)} &= -\frac{\partial \ell_j}{\partial F_\phi^{(m-1)}} = \frac{w_j(1, 1) - w_j(1, 0)e^{F_{j,\phi}^{(m-1)}}}{\left(1 + e^{F_{j,\phi}^{(m-1)}}\right)} \\ \hat{u}_{j,p}^{(m)} &= -\frac{\partial \ell_j}{\partial F_p^{(m-1)}} = \frac{w_j(1, 1) \left(1 + e^{F_{j,p}^{(m-1)}}\right) y_j - w_j(1, 1)e^{F_{j,p}^{(m-1)}}}{1 + e^{F_{j,p}^{(m-1)}}} \end{aligned} \quad (7)$$

- 258 5. For each component  $\theta = \{\phi, p\}$ :
  - 259 (a) fit  $k$  base-learners independently to the gradients:  $b_k(\hat{\mathbf{u}}_\theta^{(m)}, X_k) \Rightarrow g_k(X_k)$ ;
  - 260 (b) each fitted learner makes an estimate of the gradient,  $\hat{f}_k = g_k(X_k)$ ;
  - 261 (c) select the best-fitting base-learner  $k^* = \operatorname{argmin}_k (\hat{\mathbf{u}}_\theta^{(m)} - \hat{f}_k)^2$  and append the fitted-learner to the  
262 ensemble  $\mathcal{G}_\theta \leftarrow g_k^*$ ;
  - 263 (d) update the prediction vector:  $F_\theta^{(m)} = \nu_\theta \hat{f}_{k^*} + F_\theta^{(m-1)}$ ;
- 264 6. Estimate the empirical risk on the full data  $L(\mathbf{Y}, \mathcal{F}^{(m)})$ , or estimate the holdout-risk on a test set  
265  $L(\mathbf{Y}_{\text{test}}, \mathcal{F}_{\text{test}}^{(m)})$  s.t.  $\mathcal{F}_{\text{test}}^{(m)} = \{G_\phi^{(m)}(\mathbf{X}_{\text{test}}), G_p^{(m)}(\mathbf{X}_{\text{test}})\}$ .
- 266 7. Update  $m = m + 1$ .
- 267 8. Repeat steps 3 to 7 until  $m = m_{\text{stop}}$ .

268 The outputs of the algorithm are the fit vectors  $\mathcal{F}$  and the ensemble of fitted base-learners  $\mathcal{G}_\phi$  and  $\mathcal{G}_p$ . The  
269 estimate of survival for individual  $i$  at time  $t$  can be retrieved  $j := (i, t)$ ;  $\phi_j = \operatorname{logit}^{-1}(F_j)$ , and likewise for  
270 capture probability. For predicting  $\phi$  and  $p$  on new covariate data  $\mathbf{X}$ , we merely process the data through  
271 the ensemble of fitted base-learners and shrink by  $\nu$ , i.e.,  $F_\theta^{\text{pred}} = G_\theta(\mathbf{X}) = \nu_\theta \sum_{g_k \in \mathcal{G}_\theta} g_k(\mathbf{X})$ .

272 The three regularization parameters  $m_{\text{stop}}, \nu_\phi, \nu_p$  must be tuned by minimizing the holdout-risk averaged  
273 over many out-of-sample test sets, i.e., our estimate of the expected loss (see 2.4).

274 *2.3.2. CJSboost by Monte-Carlo approximation*

275 A second strategy to garner conditional independence of data-points  $(y_j, \mathbf{x}_j)$  and estimate the negative  
 276 gradients is to integrate over the latent state distributions  $\pi(\mathbf{z}_i | \mathbf{y}_i, \mathcal{F}_i)$  with a large sample drawn from the  
 277 posterior. A fast and simple “forward-filtering and backward-sampling” algorithm is used (Rabiner, 1989;  
 278 Murphy, 2012b), detailed in Appendix A.4. Within each boosting iteration  $m$ , we sample  $S$  sequences  
 279 of  $\mathbf{z}_i$ . Per  $s$  sequence, we estimate a separate negative-gradient, and fit base-learners to it. After fitting  
 280 all  $S$  samples, we update the prediction vectors with the best-fitting base-learners from each sequence,  
 281  $F_\theta^{(m+1)} = F_\theta^{(m)} + \nu_\theta \sum_s \hat{f}^{(s)}$ . Over  $S \times m$  draws, this is approximately equivalent to the EM algorithm. For  
 282 comparable results to CJSboost-EM, the shrinkage parameters  $\nu_{MC}$  should be set equal to  $\frac{1}{S} \nu_{EM}$ , i.e., the  
 283 contribution of any one sequence  $\mathbf{z}^{(s)}$  is small.

284 I now describe the CJSboost-MC algorithm:

- 285 1. Set parameters  $S$ ,  $m_{stop}$ ,  $\nu_\phi$ , and  $\nu_p$ .
- 286 2. Initialize  $m = 1$  and  $\mathcal{F}^{(0)}$ .
- 287 3. For  $s = 1 : S$ , do:

- 288 (a) sample latent state sequence  $\mathbf{z}_i^{(s)} \sim \pi(z | \mathbf{y}_i, \mathcal{F}_i)$  (see Appendix A.4);
- (b) estimate the negative gradients, conditional on  $\mathbf{z}_i^{(s)}$ :

$$\hat{u}_{i,t,\phi}^{(m,s)} = -\frac{\partial \ell_{i,t}}{\partial F_\phi^{(m-1)}} = \frac{\mathbb{1}[z_{i,t-1}^{(s)} = 1, z_{i,t}^{(s)} = 1] - \mathbb{1}[z_{i,t-1}^{(s)} = 1, z_{i,t}^{(s)} = 0] \cdot e^{F_{i,t,\phi}^{(m-1)}}}{1 + e^{F_{i,t,\phi}^{(m-1)}}}$$

$$\hat{u}_{i,t,p}^{(m,s)} = -\frac{\partial \ell_{i,t}}{\partial F_p^{(m-1)}} = \frac{\mathbb{1}[z_{i,t-1}^{(s)} = 1, z_{i,t}^{(s)} = 1] \left( (1 + e^{F_{i,t,p}^{(m-1)}}) y_{i,t} - e^{F_{i,t,p}^{(m-1)}} \right)}{1 + e^{F_{i,t,p}^{(m-1)}}}$$

- 289 (c) for each component  $\theta = \{\phi, p\}$ :

- 290 i. fit  $k$  base-learners independently to the gradients:  $b_k(\hat{\mathbf{u}}_\theta^{(m,s)}, X_k) \Rightarrow g_k^{(s)}(X_k)$ .
- 291 ii. each fitted learner makes an estimate of the gradient,  $\hat{f}_k^{(s)} = g_k^{(s)}(X_k)$
- 292 iii. select the best-fitting base-learner  $k^{(s)*} = \operatorname{argmin}_k (\hat{\mathbf{u}}_\theta^{(m,s)} - \hat{f}_k^{(s)})^2$  and append the fitted-learner  
 293 to the ensemble  $\mathcal{G}_\theta \leftarrow g_k^{(s)*}$ .

- 294 4. For each  $\theta = \{\phi, p\}$ : update the fit vectors, overall  $s$ :  $F_\theta^{(m)} = F_\theta^{(m-1)} + \nu_\theta \sum_s \hat{f}^{(s)}$ .

- 295 5. Estimate the empirical risk on the training data  $L(\mathbf{Y}, \mathcal{F}^{(m)})$ , or on a holdout test set  $L(\mathbf{Y}_{test}, \mathcal{F}_{test}^{(m)})$ .

- 296 6.  $m = m + 1$

- 297 7. Repeat steps 3 to 6 until  $m = m_{stop}$ .

298 Just as in the CJSboost-EM algorithm, we must tune  $\nu$  and  $m_{stop}$  through cross-validation or bootstrap-  
 299 validation (Section 2.4).

300 Notice that although the two algorithms have different specific negative-gradients and loss functions, the  
 301 empirical risk is always the negative log-likelihood of the CJS model.

## 302 2.4. Hyperparameters

303 In component-wise boosting, the three most important regularization parameters are  $m_{\text{stop}}$ ,  $\nu_\phi$ ,  $\nu_p$ . These  
304 must be tuned by some form of holdout-validation. As per Schmid et al. (2013), I suggest sampling with  
305 replacement (bootstrapping) individuals' capture histories between 50 to 100 times, training a new model on  
306 each bootstrap sample. On average, each bootstrap leaves 36.5% of the capture-histories unused in the model  
307 fitting, which can then be used to estimate a holdout-risk. Averaged over all bootstraps, this is an estimate of  
308 the generalization error. Bootstrap-validation is preferable to k-fold or leave-one-out cross-validation because  
309 it is most similar to the multiple resampling/subsampling schemes of Shao (Monte-Carlo CV and Delete-d CV;  
310 1993, 1997) which are model-selection consistent under a wider variety of conditions (e.g., sparsity, tapering).  
311 Finally, the K-bootstrap can also give us an estimate of posterior inclusion probabilities via stability-selection  
312 (Meinshausen & Bühlmann, 2010; Murphy, 2012c), which I use in section 2.7.

313 Because we can monitor the trajectory of the holdout-risk during each boosting iteration, we only need to  
314 perform one round of K-bootstrap-validation to find the optimal  $m$ . See Figure 1 for an example of monitoring  
315 the holdout-risk and estimating  $m_{\text{cv}}$ . Estimating optimal values of  $\nu_\phi$  and  $\nu_p$  is more complicated because they  
316 are continuous; in practise we must discretize the set of plausible combinations, e.g.,  $(10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}) \times$   
317  $(10^{-4}, 10^{-3}, 10^{-2}, 10^{-1})$ . Each combination requires a separate bootstrap-validation exercise. This is the  
318 most expensive step of CJSboost. See Appendix B for a suggestion on how to perform this task with only  
319 7-10 K-bootstrap-validation runs.

320 The reader should note that other multivariate boosting techniques (such as gamboostLSS; Schmid et al.,  
321 2013; Mayr et al., 2012) instead have a single fixed  $\nu$  for all parameters, and seek to optimize  $m_\theta$  per parameter  
322  $\theta$ . This is inversely related to what I propose: optimizing a global  $m_{\text{stop}}$  for both parameters, while optimizing  
323 the *ratio* of  $\nu_{\theta_1}$  to  $\nu_{\theta_2}$ . The two methods are equivalent in their outcome. In other words, making  $\nu_\theta$  smaller  
324 for component  $\theta$  is the same as decreasing  $m_\theta$  for fixed  $\nu$ , and *vice versa*. More importantly, other authors  
325 have claimed that there is little benefit in optimizing  $m$  and  $\nu$  for each component (Schmid et al., 2013).  
326 This is untrue for CJSboost, where the optimal estimate of  $\nu_\phi$  may be several orders of magnitude different  
327 than the optimal  $\nu_p$ .

328 There are theoretically many other hyperparameters, such as the base-learner parameters which control  
329 flexibility, e.g. the effective degrees-of-freedom of a spline, or the maximum tree-depth of a conditional  
330 inference tree (Hothorn et al., 2006). However, Bühlmann & Yu (2003) and Schmid & Hothorn (2008a)  
331 show that these can be safely fixed and one should instead focus on  $m_{\text{stop}}$ . A more important consideration  
332 is the *relative* flexibility of competing base-learners: multi-covariate learners and unpenalized learners have  
333 a greater freedom to (over)fit the residual variation and will be preferentially selected in the algorithm.  
334 Therefore, one should use penalties to enforce a similar effective degrees-of-freedom among all base-learners,  
335 as well as decompose higher-order interactions and non-linear curves into their constituent components. For  
336 example, if one wishes to learn about the role of covariates  $x_1$  and  $x_2$  and the possibility of an interaction  
337 between  $x_1 \times x_2$ , then one must add three PLS base-learners of equal effective-*df*: two for the main-effects

338 and a separate base-learner for their interaction. Readers should refer to the practise of “centring” in Kneib  
 339 et al. (2009) and Hofner et al. (2012).

340 *2.5. Simulation 1: MC vs EM vs AICc vs MLE*

341 The goals of this simulation were to compare estimates of survival and capture probabilities among  
 342 the two boosting algorithms (CJSboost-EM and CJSboost-MC) and benchmark them against MLEs and  
 343 AICc model-averaging. The simulated dataset was inspired by the European Dipper dataset from (Lebreton  
 344 et al., 1992). The simulated dataset included  $T = 10$  primary periods, and  $n = 300$  individuals in two  
 345 groups (male and female). Individuals’ first-capture periods ( $t_i^0$ ) were random. The true processes were  
 346 smoothly time-varying effects plus an individual covariate (sex-effect). The true data-generating processes  
 347 were:  $p(t, \text{sex}) = \text{logit}^{-1} \left( 0.5 + t \frac{\sin(t)}{17} \right) - 10 \cdot \mathbb{1}[\text{sex} = 1]$  and  $\phi(t, \text{sex}) = 0.91 - 0.01 \cdot t - 0.05 \cdot \mathbb{1}[t = 5, 6] +$   
 348  $0.05 \cdot \mathbb{1}[t = 9, 10] - 0.05 \cdot \mathbb{1}[\text{sex} = 1]$ . Figure 3 graphs the true processes.

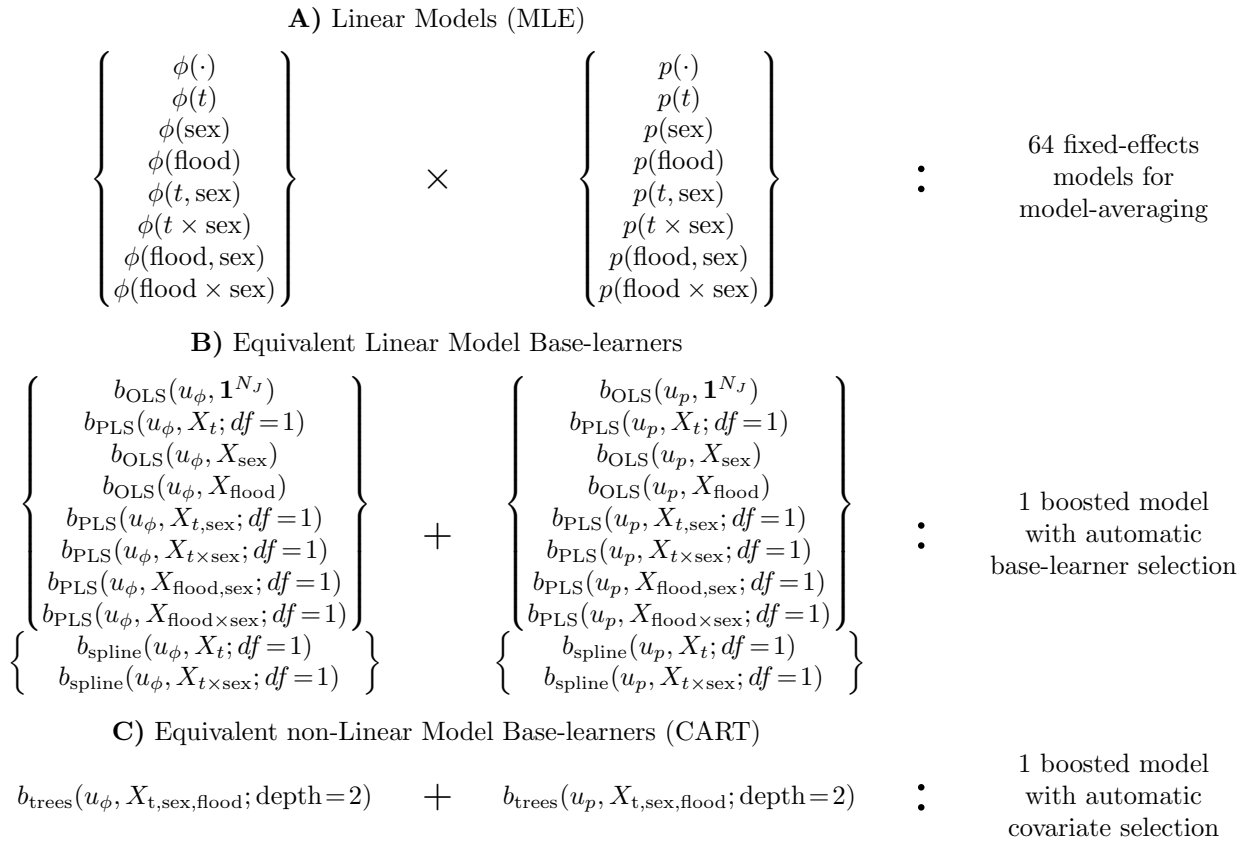


Figure 2: Different notation for multimodel inference of a Cormack-Jolly-Seber model, comparing fixed-effects model-averaging and boosting. **A)** Each fixed-effect model includes one term for  $\phi$  (left) and one for  $p$  (right).  $\theta(\cdot)$  is an intercept model;  $\theta(t)$  has different coefficients per  $T$  capture periods (with appropriate constraints on  $t=T$ );  $\theta(a, b)$  is a linear combination of covariate  $a$  and  $b$  on the logit scale;  $\theta(a \times b)$  is an interaction effect between  $a$  and  $b$  on the logit scale. **B)** Equivalent linear base-learners (Ordinary and Penalized Least Squares from `mboost`; Bühlmann & Hothorn, 2007) with penalties to constrain their effective- $df$  (ridge penalty). All terms are available in one model; selection of base-learners is by component-wise boosting. **C)** Non-linear CJS model with CART-like trees, allowing complex interactions. Selection of covariates is by the `ctree` algorithm (Hothorn et al., 2006).

349 Figure 2A shows all combinations of  $p$  and  $\phi$  parametrizations, which has 64 possible fixed-effect models  
350 for estimation by Maximum Likelihood and AICc model-averaging. The true model is best represented as  
351  $\phi(t, \text{sex})p(t, \text{sex})$ . *Flood* is a dummy categorical variable that groups the captures periods  $\{4, 5, 6\}$  (corre-  
352 sponding to a trough in either process): it simulates an analyst’s hypothesis that high flood years (in periods  
353 4,5,6) may influence dipper survival and capture probability. The MLE and AICc model-averaging analyses  
354 were conducted with Program MARK (White & Burnham, 1999) and RMark (Laake, 2013).

355 For the boosting analyses, four techniques were compared: i) linear-model CJSboost-EM (using OLS  
356 and PLS base-learners); ii) non-linear CJSboost-EM (using a CART-like base-learner called “conditional-  
357 inference trees”; Hothorn et al., 2006); iii) linear-model CJSboost-MC; and iv) non-linear CJSboost-MC.  
358 For the linear-models, the OLS and PLS base-learners included all base-learners listed in Figure 2B. See  
359 the `mboost` R package (Bühlmann & Hothorn, 2007; Hofner et al., 2012). Variable selection occurs as a  
360 consequence of the internal competition among base-learners to fit the gradient, per boosting iteration. The  
361 effective degrees-of-freedom of each base-learner were constrained with ridge penalties, as per Section 2.4.  
362 The non-linear CJSboost models had just one CART-like base-learner per  $\phi$  and  $p$ . Variable-selection and  
363 interactions are implemented internally to the `ctree` algorithm, much like a black-box.

364 All 4 models used 70-times bootstrap-validation to estimate optimal values of  $m_{\text{stop}}$ ,  $\nu_{\phi}$  and  $\nu_p$ , as per  
365 section 2.4.

### 366 2.6. Analysis: dipper example

367 Using CJSboost-EM, I reanalyzed the European Dipper dataset from (Lebreton et al., 1992). I compared  
368 the results to the MLEs of the fully-saturated model ( $\phi(t \times \text{sex})p(t \times \text{sex})$ ) as well as to AICc model-averaged  
369 estimates. The dataset has 294 individuals in  $T = 7$  capture periods. Covariates included time, sex, and  
370 flood, similar to Section 2.5. The model-building framework was the same as in Figure 2. A 70-times  
371 bootstrap-validation was used for optimizing  $m_{\text{stop}}$ ,  $\nu_{\phi}$  and  $\nu_p$ .

372 Interested readers can repeat this analysis using the online tutorial at [http://github.com/faraway1nspace/](http://github.com/faraway1nspace/HMMboost/)  
373 `HMMboost/`.

### 374 2.7. Simulation 2: high-dimensional example

375 The final simulation addressed the issue of high-dimensionality and the ability of CJSboost (EM) to find  
376 a sparse set of important covariates out of many spurious covariates. This is a variant of the “small n big p”  
377 problem often studied in machine learning. However, this challenge is extraordinarily difficult for capture-  
378 recapture analysis, because one must consider all combinations of covariates for different parameters ( $\phi, p$ ).  
379 In this section, I simulated 21 multi-colinear, individual-level covariates (18 continuous, three discretized)  
380 drawn from a multivariate Gaussian with marginal variances of 1. The general model can be represented as:

$$\text{logit}(\theta_{i,t}) = \beta_{\theta,0} + \underbrace{\sum_{k=1}^{21} X_{i,k} \beta_{\theta,k}}_{\text{individual effects}} + \underbrace{\sum_{\tau=2}^T \beta_{\theta,\tau} \mathbb{1}[\tau=t]}_{\text{capture period effect}}$$

381 The intercepts were drawn randomly from  $\beta_{p,0} \sim U(0.4, 0.6)$  and  $\beta_{\phi,0} \sim U(0.55, 0.8)$ . The true models  
382 were deliberately *sparse*, such that only three covariates' coefficients ( $\beta_{\theta}^*$ ) were non-zero. For continuous  
383 variables,  $\beta_{\theta}^*$  had a norm of 1 (on the logit scale), while the categorical-variables had norms of 3, resulting  
384 in individual marginal effects spanning 0.8–0.9 probability-units. Time-as-a-categorical-variable was also  
385 included as a possible influential covariate. The number of individuals varied randomly from  $n = 200:300$ , in  
386  $T = 10$  capture periods. The simulation was repeated 30 times, each time with new covariates and coefficients.

387 Such a model-averaging exercise cannot be performed in MARK, because there are more than 4 trillion  
388 different fixed-effects models (excluding two-way interactions or higher). Furthermore, the AIC is known to  
389 do poorly when the simulated true model is sparse by design (Burnham, 2004).

390 For each simulated dataset, the boosting analyses used 23 different PLS base-learners ( $df = 2$ ) for all  
391 continuous and categorical covariates, and included capture-period  $t$  as a categorical variable (a.k.a, the  $\theta(t)$   
392 model). A 70-times bootstrap-validation was performed to optimize  $m_{\text{stop}}$ ,  $\nu_p$ , and  $\nu_{\phi}$ . After optimization,  
393 the performance of the fitted models were assessed by calculating 2 point-wise statistics between the true  
394 (simulated) processes and the estimates of  $\text{logit}(\phi)$  and  $\text{logit}(p)$ : i) Pearson correlation  $\rho(F_{\theta}^{(\text{true})}, \hat{F}_{\theta})$ ; and ii)  
395 the slope between  $s(F_{\theta}^{(\text{true})}, \hat{F}_{\theta})$  from a simple linear regression, whereby  $s = 1$  suggests that the estimates  
396 are unbiased.  $\hat{\rho}_{\theta}$  is a measure of the precision of the linear relationship between the true and fitted values,  
397 while  $\hat{s}_{\theta}$  can be likened to angular bias.

398 An extra topic explored in the online tutorial, but not in this paper, was the performance of CART-like  
399 trees (see <http://github.com/faraway1nspac/HMMboost/>).

400 In addition to studying the precision and bias of estimates, I also demonstrate the usefulness of inclusion  
401 probabilities (the probability that a covariate is selected in the model) to infer the importance of covariates. I  
402 used the technique of stability selection from Meinshausen & Bühlmann (2010), integrated within the 70-times  
403 bootstrap-validation. Stability selection probabilities  $\hat{S}$  are estimated by scoring whether a  $k^{\text{th}}$  covariate  $X_k$   
404 is selected in a  $b$  bootstrap before  $m$  iterations,  $\hat{S}_{\theta,k}^{(m)} = \frac{1}{70} \sum_{b=1}^{70} \mathbb{1}[X_k \in \mathcal{G}_{\theta}^{(m,b)}]$ ;  $\hat{S}_{\theta,k}^{(m)}$  is evaluated per  $m$ , per  
405 covariate  $X_k$  and per parameter  $\theta \in \{\phi, p\}$ . The average over all (reasonable) regularization hyperparame-  
406 ters yields a Frequentist approximation to posterior inclusion probabilities,  $\pi(i_{\theta,k}|\mathcal{D}) \approx \frac{1}{m_{\text{max}}} \sum_{m=1}^{m_{\text{max}}} S_{\theta,k}^{(m)}$   
407 (Murphy, 2012c). Ideally, influential covariates should have very high inclusion probabilities ( $\gg 0.5$ , and  
408 perhaps close to 1). Such posterior probabilities are an important means of inference about the covariates,  
409 and are more intuitive than other familiar tools for inference, like 95%CI (Hoekstra et al., 2014; Morey et al.,  
410 2016). Also, the *stability paths* (Figure 8) can be a valuable graphical tool for interpreting the importance of  
411 covariates (Meinshausen & Bühlmann, 2010).

412 Stability selection can also serve in a second-stage of “hard-thresholding” to find a sparse set of truly  
413 influential covariates (Bach, 2008; Meinshausen & Bühlmann, 2010). One picks an inclusion probability  
414 threshold between 0.5–0.99, and discards non-influential covariates below this threshold. One can proceed  
415 to “debias” the coefficients by running a final boosting model using only the selected covariates (Murphy,  
416 2012c) and setting  $m \rightarrow \infty$ . Choosing an appropriate threshold is a classic trade-off between Type I errors



417 and Power: a high threshold  $\approx 1$  should correctly reject the non-influential covariates (low False Discovery  
418 Rate) but may wrongly reject some of the truly influential covariates (high False Rejection Rate); a low  
419 threshold  $< 0.5$  will result in a higher False Discovery Rate but low False Rejection Rate. Ideally, there  
420 should be a wide range of thresholds between 0.5-1 where both the FDR and FRR are close to zero.

421 When the FDR and FRR are zero, a procedure is called “model-selection consistent”: it can correctly  
422 shrink the coefficients of non-influential covariates to zero. It is also an “oracle” if it can accurately estimate  
423 coefficients as if the true model was known in advance (Zou, 2006). The Lasso, Ridge-regression, and Boosting  
424 do *not* have these properties (Zou, 2006; Bach, 2008; Bühlmann & Hothorn, 2010): there is a pernicious trade-  
425 off between predictive-performance and model selection consistency (Zou, 2006; Meinshausen & Bühlmann,  
426 2006; Murphy, 2012c) which has to do with one’s values (Vrieze, 2012). The AIC is also not model-selection  
427 consistent (Shao, 1997; Vrieze, 2012). Instead, the AIC and Boosting are motivated by good prediction  
428 performance and minimizing the expected loss, rather than the belief in a sparse true model. Many authors  
429 laud this latter perspective, and declare sparsity to be a purely human construct that is irrelevant to natural  
430 phenomena (Burnham, 2004; Vrieze, 2012). Philosophical notions aside, there may be a strong practical  
431 imperative in capture-recapture to favour sparser solutions than what AIC or boosting can provide, as  
432 we demonstrate with the stability paths. Towards this goal, stability selection and inclusion probabilities  
433 can transform an  $\ell_1$  regularizer into a model-selection consistent procedure (Bach, 2008; Meinshausen &  
434 Bühlmann, 2010). Further debiasing can give it oracle properties. Such a multi-stage procedure is no longer  
435 strictly about prediction; rather, it considers regularization as an intermediary step towards an ultimate goal  
436 to recover a true sparse model.

437 One caveat to using stability selection for CJSboost is that base-learners must have equal flexibility/degrees-  
438 of-freedom; otherwise, the more complex base-learners (and their constituent covariates) will have a greater  
439 probability of being selected (Kneib et al., 2009). See Section 2.4.

440 A final note on debiasing and convexity of the loss function: after thresholding, the final model may not  
441 have a unique MLE, such as as the  $\phi(t)p(t)$  model. In such cases, one must impose constraints (such as  
442  $\phi_T = \phi_T$ ) before attempting to debias the results and run the gradient descent until convergence  $m \rightarrow \infty$ .  
443 Regularized CJSboosting does not have this problem because of early-stopping and model-selection.

### 444 3. Results

#### 445 3.1. Simulation 1: EM vs MC vs AICc vs MLE

446 Figure 3 compares the estimates from CJSboost-EM versus AICc model-averaging and MLEs from the  
447 full-model  $\phi(t \times \text{sex})p(t \times \text{sex})$ , as well as the true processes. Figure 4 does the same for the CJSboost-MC.  
448 The results can be summarized as follows:

- 449 i) The Expectation-Maximization algorithm and the Monte-Carlo algorithm yielded approximately the  
450 same estimates for the linear models ( $b_{\text{PLS}}$ ), but slightly different results for the non-linear CART-like  
451 base-learners ( $b_{\text{Trees}}$ ).

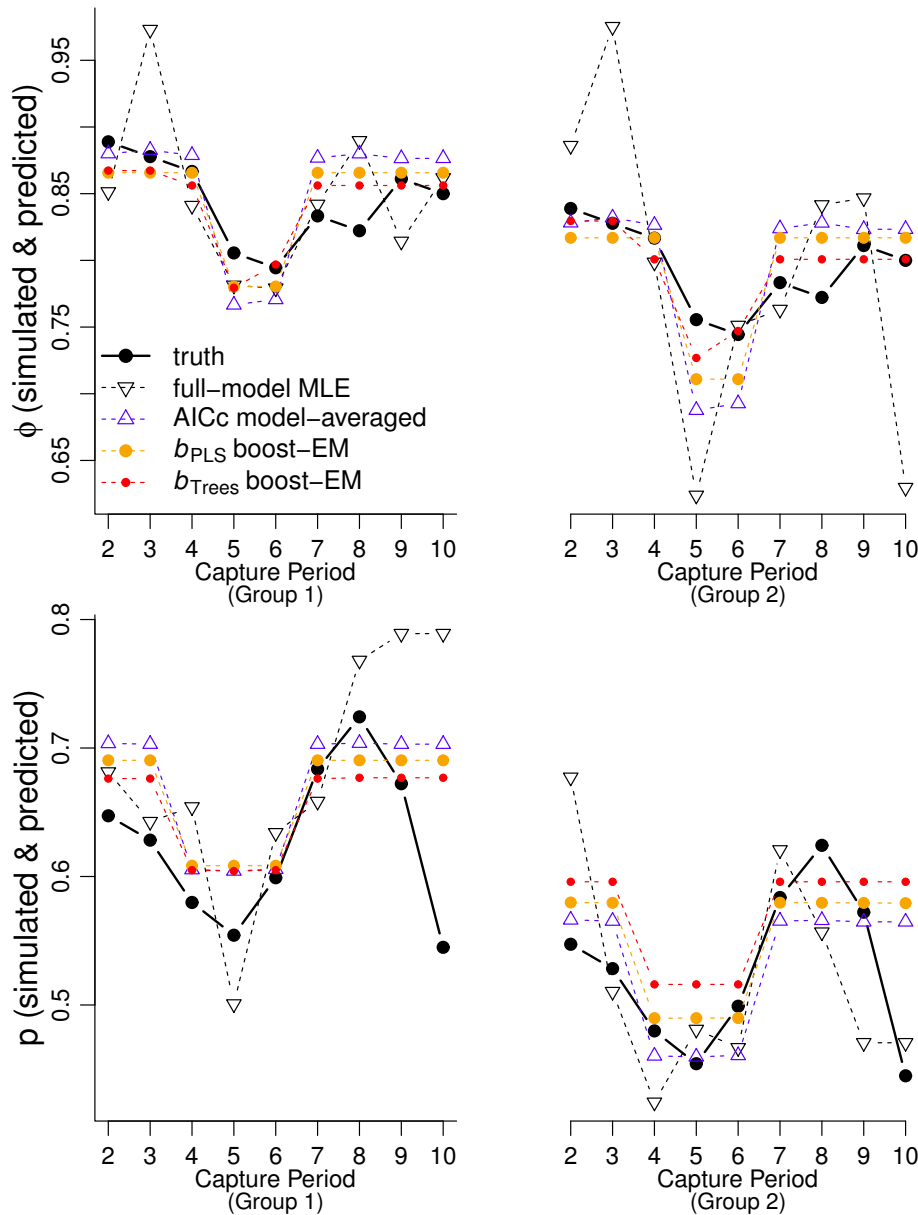


Figure 3: Simulation 1, demonstrating the CJSboost estimates from the Expectation-Maximization technique. A comparison of capture probability estimates  $\hat{p}(t \times \text{sex})$  and survival estimates  $\hat{\phi}(t \times \text{sex})$  from models composed of linear base-learners (OLS and PLS; in orange) and non-linear base-learners (CART-like trees; in red), as well AICc model-averaging (blue) and MLE (dashed black).

- 452 ii) None of the four methods (MLE, AICc,  $b_{\text{PLS}}$ -boost or  $b_{\text{trees}}$ -boost) did a convincing job of approximating  
 453 the true underlying processes (for both  $\phi$  and  $p$ ), although each model did uncover some aspect of the  
 454 true processes.
- 455 iii) The similarities between the three predictive methods (AIC,  $b_{\text{PLS}}$ -boost,  $b_{\text{trees}}$ -boost) were thus:  
 456 (a) all three methods showed the same pattern for both for  $\phi$  and  $p$ : low values during the flood periods  
 457 ( $t=4, 5, 6$ ), and a moderate sex effect (group 1 had higher values than group 2);  
 458 (b) the  $b_{\text{PLS}}$ -boost model was most similar to AICc model-averaging;

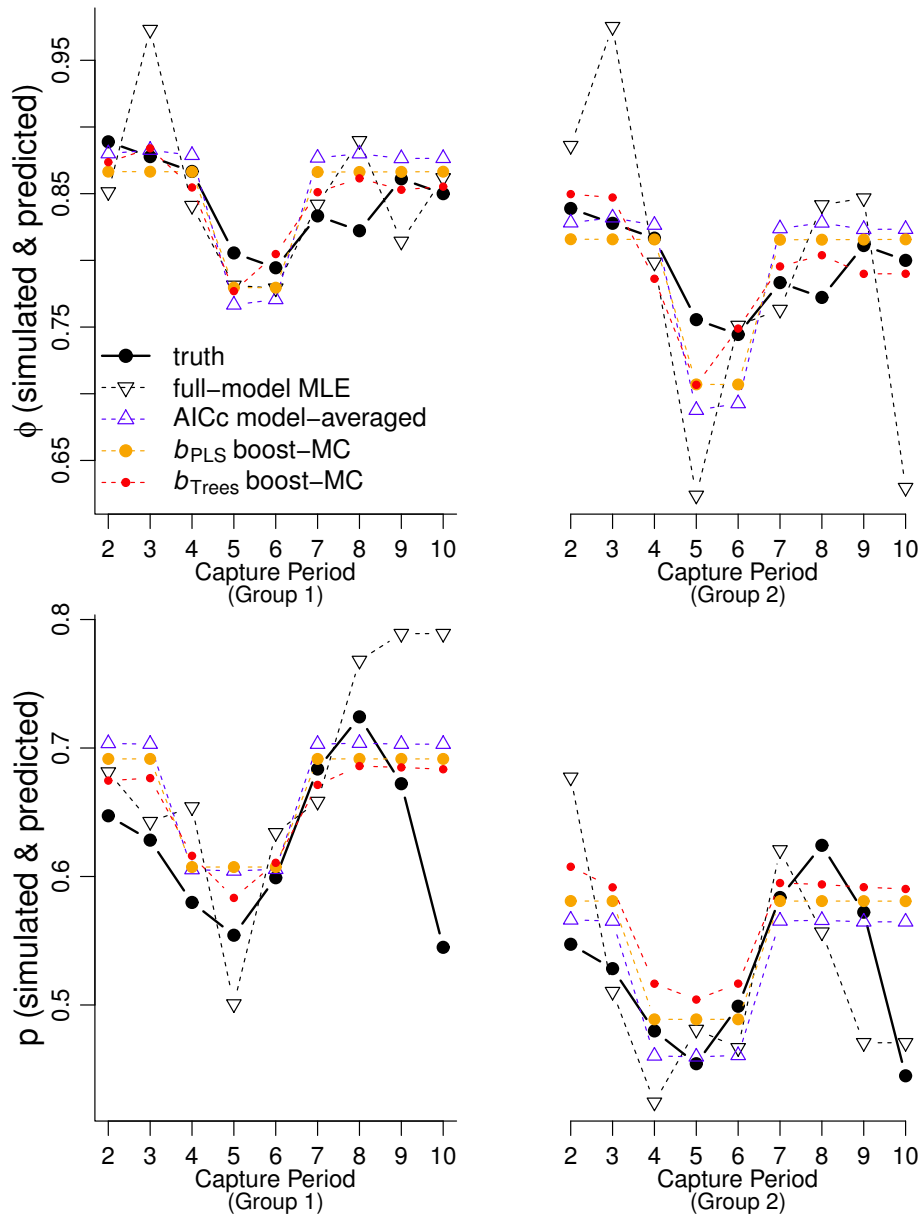


Figure 4: Simulation 1, demonstrating CJSboost estimates from the Monte-Carlo approximation technique. A comparison of capture probability estimates  $\hat{p}(t \times \text{sex})$  and survival estimates  $\hat{\phi}(t \times \text{sex})$  from models composed of linear base-learners (OLS and PLS; in orange) and non-linear base-learners (CART-like trees; in red), as well AICc model-averaging (blue) and MLE (dashed black).

- 459 (c) the estimates from both boosted models were *shrunk to the mean* relative to model-averaged esti-  
 460 mates; i.e., high model-averaged estimates were generally greater than the boosted estimates, and  
 461 low model-averaged estimates were generally lower than the boosted estimates;  
 462 (d) the non-linear  $b_{trees}$  estimates showed more shrinkage to the mean than the linear  $b_{PLS}$  estimates;  
 463 iv) The MLEs of the full-model  $\hat{\phi}(t \times \text{sex})\hat{p}(t \times \text{sex})$  showed the most extremes values.

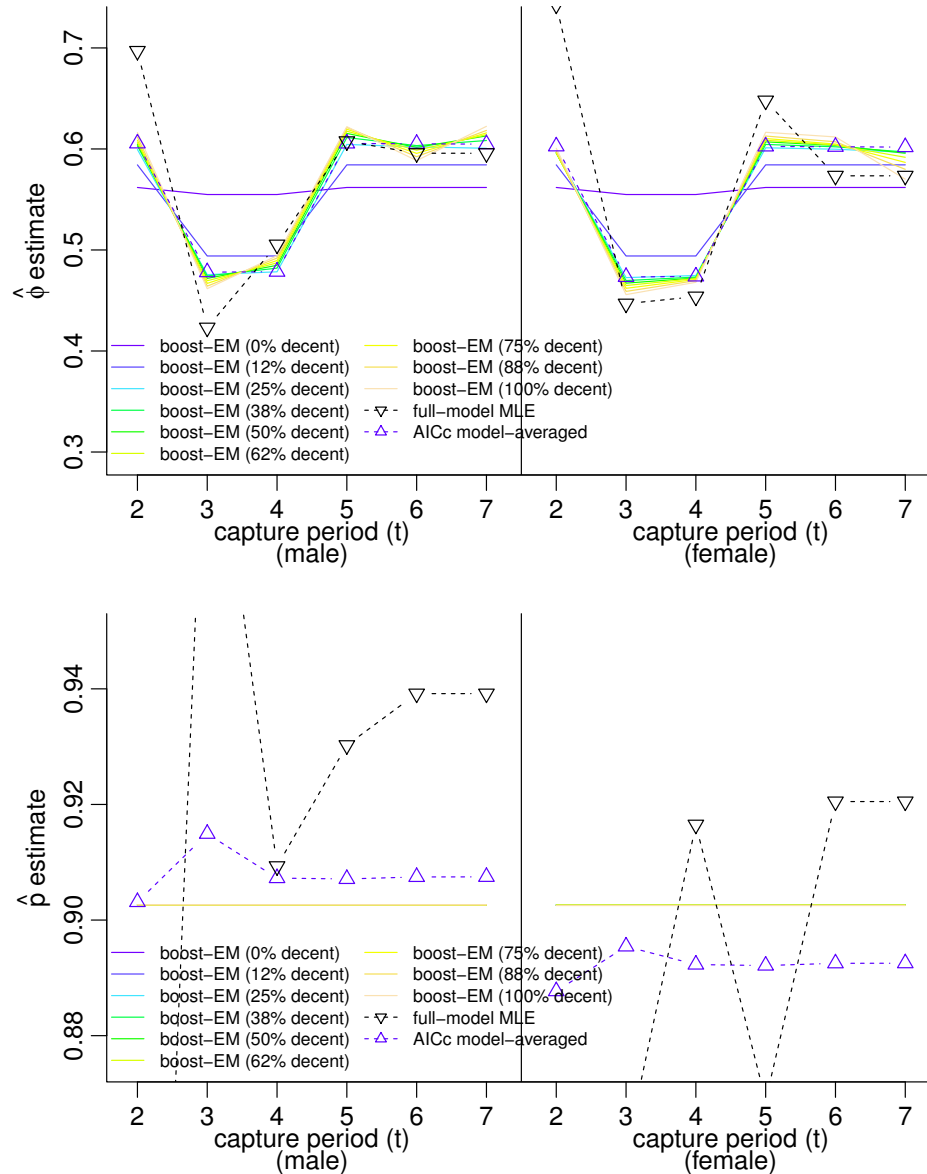


Figure 5: Dipper analysis of survival  $\phi$  and capture probability  $p$  by CJSboost-EM using least-squares base-learners, plus comparison with AICc model-averaging and MLE ( $\hat{\phi}(t \times \text{sex})\hat{p}(t \times \text{sex})$ ). The regularization pathway of the estimates is visualized with the spectrum-coloured lines, starting at the intercept-only model (0% descent) and growing more complex as the gradient descent algorithm proceeds. The final estimates are achieved at 100% of the descent, when the boosting iteration  $m_{CV}$  is reached.

### 464 3.2. Results: Dipper example

465 This section shows the reanalysis of the European Dipper dataset from Lebreton et al. (1992). Figure  
 466 1 shows an example of the gradient descent of the empirical risk and the holdout-risk from the 70-times  
 467 bootstrap-validation used to estimate the optimal  $m_{\text{stop}}$ . Comparisons were between the linear  $b_{\text{PLS}}$ -boost-  
 468 EM model (Figure 5) and the non-linear  $b_{\text{Trees}}$ -boost-EM model (Figure 6), as well as AICc model-averaging  
 469 and MLEs from the full-model  $\phi(t \times \text{sex})p(t \times \text{sex})$ . Both Figures also show the “regularization pathway” of

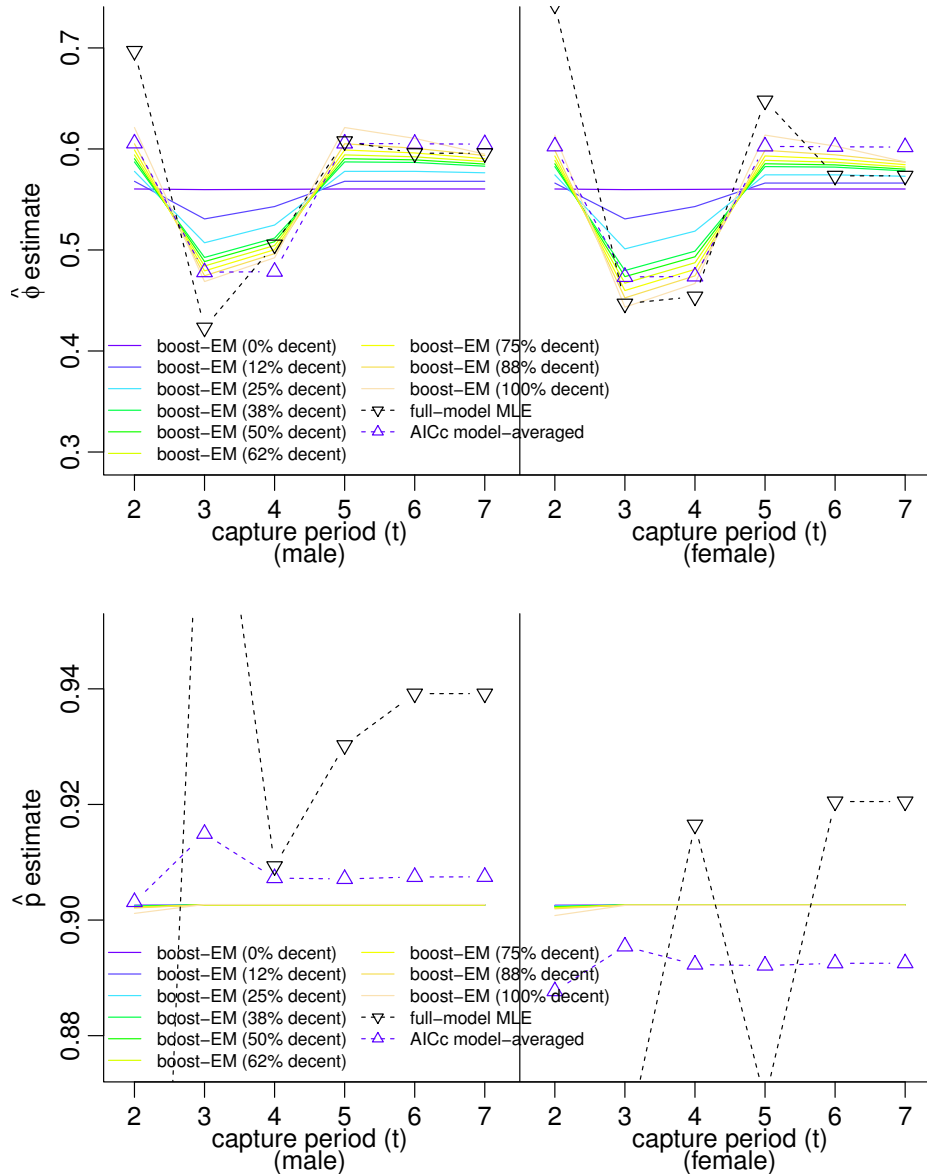


Figure 6: Dipper analysis of survival  $\phi$  and capture probability  $p$  by CJSboost-EM using non-linear base-learners (CART-like trees), plus comparison with AICc model-averaging and MLE ( $\hat{\phi}(t \times \text{sex}) \hat{p}(t \times \text{sex})$ ). The regularization pathway of the estimates is visualized with the spectrum-coloured lines, starting at the intercept-only model (0% decent) and growing more complex as the gradient descent algorithm proceeds. The final estimates are achieved at 100% of the descent, when the boosting iteration  $m_{CV}$  is reached.

470 their respective boosted model: the movement of the estimates from their initial uniform intercept model (at  
 471  $m = 0$ ) to their final values at  $m = m_{CV}$ , stratified by the percentage of the total reduction in the empirical  
 472 risk. The results can be summarized thus:

- 473 i) For both survival  $\phi$  and capture probability  $p$ , all three predictive methods (AICc,  $b_{\text{PLS-boost}}$  or  $b_{\text{trees-boost}}$ )  
 474 were much more similar to each other than to the MLEs from the full-model.  
 475 ii) For survival, all three predictive methods yielded the same estimates: a survival probability of 0.48-0.5

476 during the flood years ( $t=3, 4$ ) and no sex-effect.  
477 iii) For capture probability, the model-average estimates suggested a slight sex effect of about 1.5 probability  
478 units, whereas both boosted models shrunk the capture-probability to a constant; in contrast, the MLEs  
479 varied wildly.  
480 iv) Regarding the regularization pathways, the linear  $b_{\text{PLS}}$ -boosted estimates converged very quickly (within  
481 25% of the gradient decent) to their final estimates; whereas the movement of the non-linear  $b_{\text{Trees}}$ -  
482 boosted estimates moved much more gradually.

483 *3.3. Simulation 2: high-dimensional example*

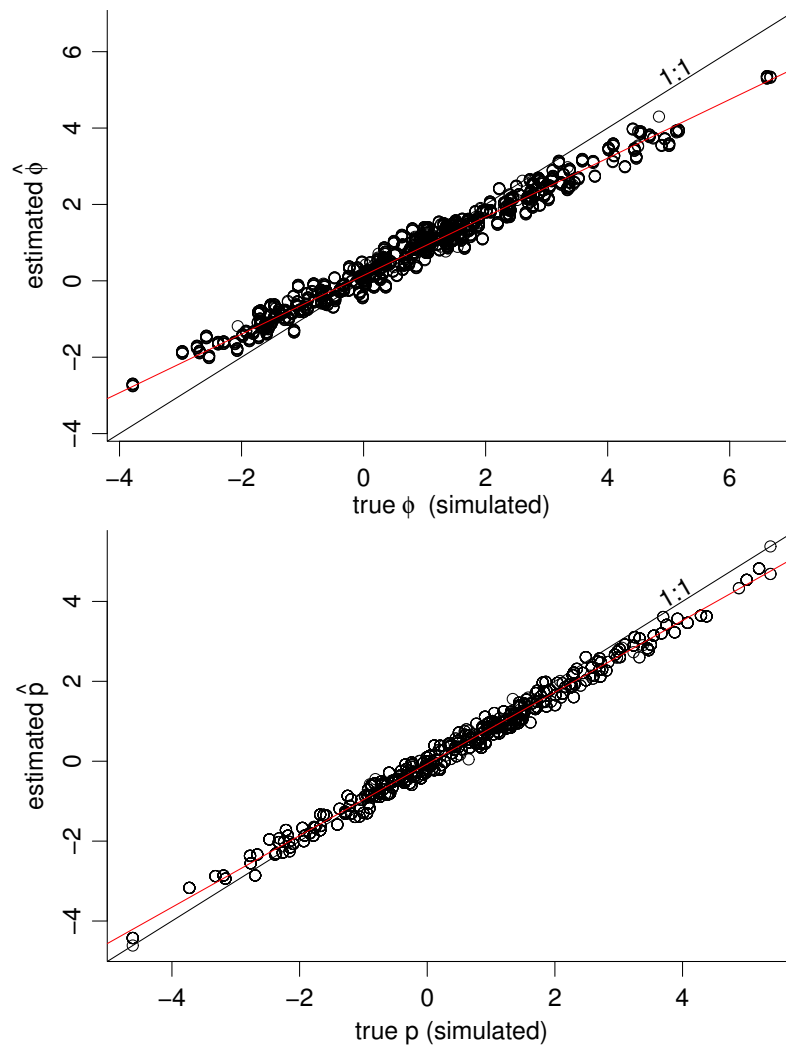


Figure 7: Comparing true vs estimated survival  $\phi_{i,t}$  and capture probability  $p_{i,t}$  for individuals  $i$  at capture-period  $t$ . Boosted estimates incur some downward bias (evident in the difference between the 1:1 line and the estimates' red trend-line) due to shrinkage of coefficients to the intercept-only model.

484 Over the 30 simulations, the Pearson correlation between the true and estimated survival had the following  
485 descriptive statistics: mean of 0.979, minimum of 0.955,  $Q_{0.05}$  of 0.959, and a maximum of 0.997. For capture

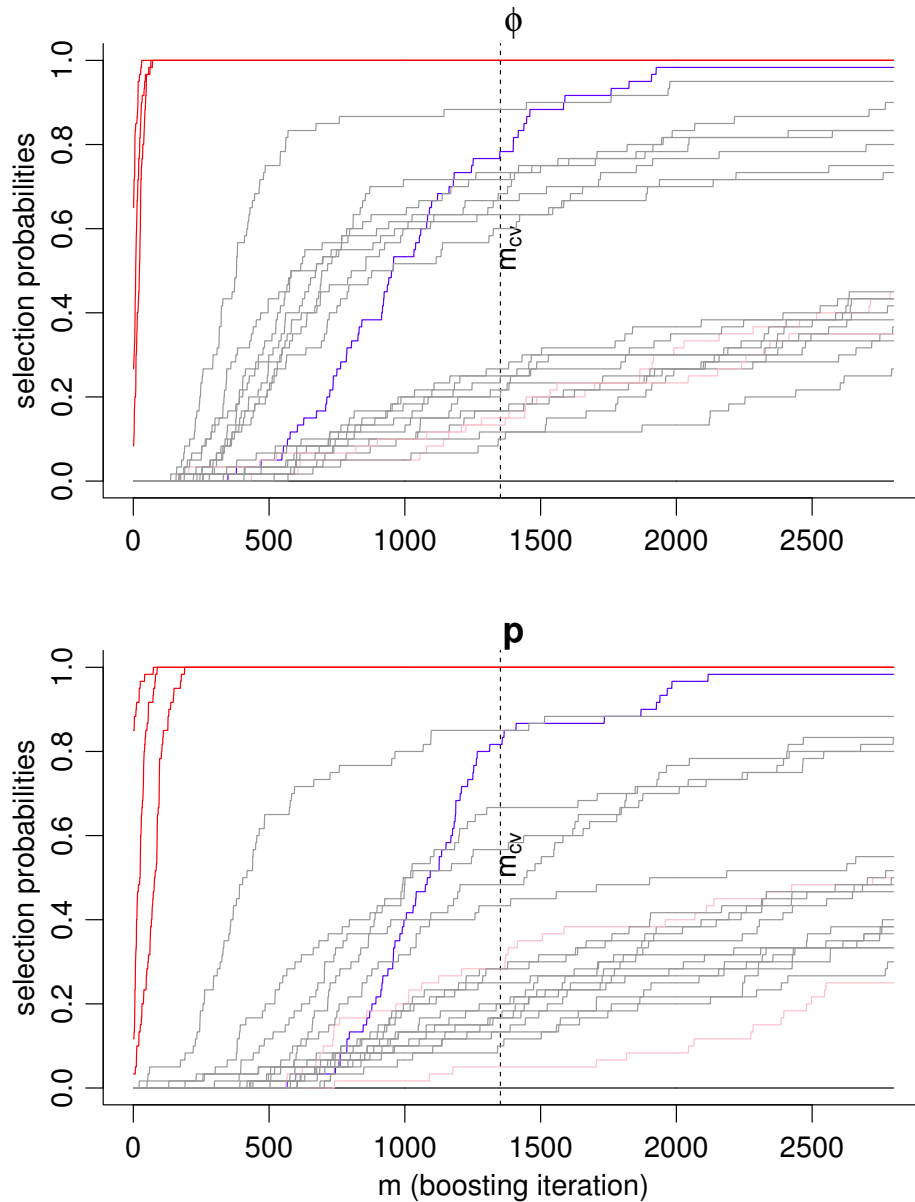


Figure 8: Demonstration of stability selection probabilities for one high-dimensional simulation. As the boosting iteration ( $m$ ) gets large, regularization gets weaker, and all covariates have a higher selection probability  $S$  (estimated from a bootstrap). Lines in **red** are truly influential covariates. Lines in **gray** are non-influential covariates. Lines in **pink** are not-influential for  $\theta$ , but are influential in the other parameter  $-\theta$ . Lines in **blue** represent the time-as-a-categorical-variable base-learner, a.k.a  $\theta(t)$ , which in this simulation was non-influential.

486 probability, the same statistics were: 0.961, 0.708, 0.911 and 0.998. The slope statistic, a measure of bias  
 487 between estimated and true survival, had the following statistics: mean of 0.778, minimum of 0.618,  $Q_{0.05}$   
 488 of 0.647, and maximum of 1.018. For capture probabilities, these slope statistics were: 0.782, 0.404, 0.542,  
 489 0.962. Figure 7 shows the results of one example simulation to demonstrate the high-precision and slight-bias  
 490 that is characteristic of boosting algorithms and other  $\ell_1$  regularizers.

491 Regarding the stability selection results, Figure 8 shows an example of the stability paths over  $m$  (for the



492 same simulations shown in Figure 7). Readers can view an online animated GIF which shows the stability  
493 paths for all 30 simulations, at <http://github.com/farawayinspace/HMMboost/> and in the Supplementary  
494 Material. The results can be summarized as:

- 495 i) The stability paths of the truly influential covariates (in *red*, Figure 8) were visually very different from  
496 the rest of the non-influential covariates (*grey*). In particular, the truly influential covariates reached  
497 high stability selection probabilities  $S$  for small values of  $m$ . For  $\phi$ , they reach  $S_{\phi}^{(m_{CV})} = 1$  by the optimal  
498  $m_{CV}$  in all simulations; while for  $p$ , 94.6% of the covariates reached  $S_p^{(m_{CV})} = 1$  by  $m_{CV}$ . On average,  
499 their posterior inclusion probabilities (over all  $m$ ) were 0.98 and 0.96 for  $\phi$  and  $p$ , respectively.
- 500 ii) For the non-influential covariates, the stability selection probabilities at  $m_{CV}$  were low,  $S^{(m_{CV})} \lesssim 0.5$ ,  
501 and rarely achieved a  $S^{(m_{CV})} > 0.8$  by  $m_{CV}$ . Only 1.2% of such covariates achieved  $S^{(m_{CV})} \geq 0.95$  by  
502  $m_{CV}$ , for both  $\phi$  and  $p$ . On average, their inclusion probabilities were 0.32 for  $\phi$  and 0.38 for  $p$ .
- 503 iii) The stability path of the time-as-a-categorical-variable (a.k.a  $\theta(t)$ , in *blue*, Figure 8) showed a greater  
504 tendency for inclusion and achieved high stability selection probabilities, particularly for  $p$ . For  $p$ , it  
505 achieved  $S_p^{(m_{CV})} \geq 0.95$  by  $m_{CV}$  in 60% of simulations (in which it was not truly influential). Its inclusion  
506 probabilities were 0.49 for  $\phi$  and 0.75 for  $p$ , averaged over all simulations. This has important implications  
507 for model-selection consistency (or lack thereof). This may explain the anecdotal experience that, to have  
508 good-fitting capture-recapture models, one must usually incorporate time-varying capture-probabilities.
- 509 iv) The stability paths of covariates which were important in one parameter (like  $\phi$ ) but unimportant in the  
510 other parameter (like  $p$ ) seemed to achieve higher inclusions probabilities (in *pink*, Figure 8), more so  
511 than the other non-influential covariates in grey. For  $p$ , such covariates achieved  $S_p^{(m_{CV})} \geq 0.95$  in 10%  
512 of simulations, and in 3.3% of simulations for  $\phi$ . This suggests an underlying structural correlation and  
513 may have implications for model-selection consistency.

514 Table 1 shows the coefficients of a prediction-optimal CJSboost model for one simulation (same as Figures  
515 7 and 8). As expected, the regularized regression coefficients placed highest weight on the 6 truly influential  
516 covariates, albeit with a downward bias that is characteristic of  $\ell_1$  regularization (the true values were  $\|\beta_k\| =$   
517 1). The model shrunk the remaining non-influential coefficients to low values, but not to zero, incurring a  
518 False Discovery Rate of 0.34. Table 1 also demonstrates the effects of coefficient hard-thresholding using the  
519 posterior inclusion probabilities estimated in Figure 8. At higher thresholds (0.80-0.95), the model succeeds  
520 in having a FDR and FRR of zero, as well as accurate unbiased estimates of the coefficients (seemingly an  
521 “oracle”, for this one simulation). The optimal threshold seems to be in the of 0.80-0.95, similar to the  
522 threshold suggested by Bach (2008) and Meinshausen & Bühlmann (2010). After “debiasing” (Murphy,  
523 2012c; here, meaning running  $m \rightarrow \infty$  after hard-thresholding), the CJSboost estimates become nearly equal  
524 to the oracle MLEs (a benchmark model run with 100% foresight about the true model). Thresholding at low  
525 values ( $< 0.8$ ) and debiasing added too much weight on some non-influential covariates (i.e., no shrinkage),  
526 whereas thresholding at extremely high values ( $> 0.95$ ) incurred a False Rejection.

Table 1: Estimates of coefficients from CJSboost, for one high-dimensional model-selection problem, under different degrees of hard-thresholding

Parameter	Survival $\Phi$									MLE Oracle <sup>§</sup>	SE Oracle
	Prediction Optimal <sup>†</sup>	0.55	0.65	Inclusion Probability Threshold <sup>‡</sup>					0.95		
				0.75	0.8	0.85	0.9				
$\hat{\beta}_\phi$ (time:1)	-0.002	-0.01	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (time:2)	-0.041	-0.238	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (time:3)	-0.036	-0.271	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (time:4)	-0.026	-0.285	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (time:5)	0.017	0.205	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (time:6)	0.006	-0.005	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (time:7)	0.015	0.124	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (time:8)	0.022	0.196	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (time:9)	0.025	0.264	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (time:10)	-0.001	-0.091	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (a)	-0.083	-0.173	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (b)	0.828	0.982	1.064	1.045	<b>1.067</b>	<b>1.067</b>	<b>1.067</b>	<b>1.067</b>	1.074	1.068	0.143
$\hat{\beta}_\phi$ (c)	-0.021	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (d)	-0.761	-0.93	-0.991	-0.983	<b>-0.965</b>	<b>-0.965</b>	<b>-0.965</b>	<b>-0.965</b>	-0.919	-0.967	0.123
$\hat{\beta}_\phi$ (e)	0.175	0.262	0.288	0.303	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (f)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (g)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (h)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (i)	-0.051	-0.107	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (j)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (k)	-0.717	-0.838	-0.975	-0.968	<b>-0.953</b>	<b>-0.953</b>	<b>-0.953</b>	<b>-0.953</b>	-0.868	-0.955	0.119
$\hat{\beta}_\phi$ (l)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (m)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (n)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (o)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (p)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (q)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (r)	-0.048	-0.151	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (s:1)	-0.034	-0.109	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (s:2)	0.028	0.093	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (t:1)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (t:2)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (u:1)	-0.061	-0.165	0	0	0	0	0	0	0	0	0
$\hat{\beta}_\phi$ (u:2)	0.059	0.166	0	0	0	0	0	0	0	0	0
	Capture Probability $p$										
$\hat{\beta}_p$ (time:1)	0	0.002	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (time:2)	0	0.266	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (time:3)	0	-0.23	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (time:4)	0	-0.041	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (time:5)	0	-0.098	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (time:6)	0	0.159	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (time:7)	0	-0.04	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (time:8)	0	0.123	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (time:9)	0	-0.056	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (time:10)	0	-0.062	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (a)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (b)	0.942	1.129	1.149	1.184	<b>1.176</b>	<b>1.176</b>	<b>1.176</b>	<b>1.176</b>	0.846	1.178	0.144
$\hat{\beta}_p$ (c)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (d)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (e)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (f)	-0.933	-1.142	-1.181	-1.189	<b>-1.186</b>	<b>-1.186</b>	<b>-1.186</b>	<b>-1.186</b>	-0.856	-1.189	0.135
$\hat{\beta}_p$ (g)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (h)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (i)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (j)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (k)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (l)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (m)	0.042	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (n)	0.01	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (o)	0.81	0.993	1.033	1.047	<b>1.059</b>	<b>1.059</b>	<b>1.059</b>	<b>1.059</b>	0	1.061	0.124
$\hat{\beta}_p$ (p)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (q)	-0.027	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (r)	-0.063	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (s:1)	-0.15	-0.202	-0.243	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (s:2)	0.116	0.161	0.197	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (t:1)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (t:2)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (u:1)	0	0	0	0	0	0	0	0	0	0	0
$\hat{\beta}_p$ (u:2)	0	0	0	0	0	0	0	0	0	0	0
<b>False Discovery Rate:</b>	0.342	0.237	0.053	0.026	0	0	0	0	0	0	0
<b>False Rejection Rate:</b>	0	0	0	0	0	0	0	0	0.167	0	0

**Bold coefficients** show oracle-properties.

Covariates  $a-r$  are continuous; covariates  $s-u$  are categorical;  $\beta$ (time:t) is equivalent to a  $\theta(t)$  sub-model.

<sup>†</sup> CJSboost-EM model with  $m_{\text{stop}}$  tuned by bootstrap-validation.

<sup>‡</sup> Debiased CJSboost-EM model (unregularized;  $m \rightarrow \infty$ ) after discarding covariates with inclusion probabilities below a threshold.

<sup>§</sup> MLEs when the true model is known in advance.

## 527 4. Discussion

528 This study presents a boosted ensemble method for the Cormack-Jolly-Seber capture-recapture model,  
529 called CJSboost. I compared its estimates to AICc model-averaging. While univariate boosting is well-known  
530 in applied ecology (Elith et al., 2008; Hothorn et al., 2010; Tyne et al., 2015), the naive application of boosting  
531 for capture-recapture was not possible because of the serially-dependent nature of capture-histories. In  
532 response to this challenge, this paper presents two modifications to the Component-wise Boosting procedure,  
533 one based on Expectation-Maximization (first suggested in the appendix of Ward et al., 2009) and another  
534 based on Monte-Carlo imputation of HMM latent states. Both lead to equivalent inferences (up to an  
535 approximation error) and serve to validate each other. Code and a tutorial are available on the Github site  
536 <http://github.com/farawayinspace/HMMboost>. The framework can be easily extended to other capture-  
537 recapture systems, thereby introducing new machine-learning techniques to capture-recapture practitioners,  
538 such as CART-like trees, splines and kernels.

539 The motivation for boosted capture-recapture models are many:

- 540 1. automatic variable selection and step-wise multimodel inference (without the sometimes-impossible task  
541 of fitting all possible fixed-effects models, as in AIC-based model averaging);
- 542 2. regularization and sparse estimation, which deflate the influence of unimportant covariates;
- 543 3. shrinkage of estimates away from extreme values and inadmissible values (e.g.,  $\phi = 1$ );
- 544 4. a smoother way to address parameter non-identifiability issues, via regularization and step-wise esti-  
545 mation, rather than arbitrary constraints (e.g., fixing  $\phi_T = \phi_{T-1}$ );
- 546 5. highly extensible (see the wide variety of base-learners available under the `mboost` package, Bühlmann  
547 & Hothorn, 2007; Hofner et al., 2012);
- 548 6. inference based on predictive performance.

549 Through simulation and an analysis of the European Dipper dataset (Lebreton et al., 1992), this study  
550 is primarily concerned with comparisons of CJSboost to AICc model-averaging. This is not because of  
551 theoretical connections between the two (although some do exist); rather, AIC model-selection and model-  
552 averaging are the incumbent multimodel inference techniques in capture-recapture practise. It is therefore  
553 very reassuring that estimates from CJSboost and AICc model-averaging are qualitatively comparable, re-  
554 vealing strikingly similar patterns. This was apparent among simple least-squares base-learners as well as  
555 purely-algorithmic base-learners like CART. One distinction was that the CJSboost models were slightly more  
556 conservative and had more shrinkage on coefficients. This is desirable, especially during the current crisis  
557 of reproducibility (Simmons et al., 2011; Yaffe, 2015), because the AIC is thought to be overly permissive  
558 (Shao, 1993, 1997; Burnham, 2004; Vrieze, 2012; Hooten & Hobbs, 2015).

559 Secondly, the AIC serves as a useful conceptual bridge for introducing practitioners to the notion of  
560 regularization and predictive performance. For instance, the AIC is itself a specific type of regularized  
561 objective function (fixed-penalty  $\ell_0$  regularizer) nested within a more general class of regularizers, within

562 which Component-wise Boosting is generally considered a  $\ell_1$  regularizer (Efron et al., 2004; Bühlmann &  
563 Hothorn, 2007). The AIC also has a cross-validation interpretation (Stone, 1977; Shao, 1993, 1997). There-  
564 fore, capture-recapture practitioners, who are already (perhaps unwittingly) using predictive-performance  
565 and regularization, should expand their concept of “model parsimony” and multi-model inference to include  
566 boosting. There has been a call for ecologists to embrace algorithmic means of inference (Oppel et al., 2009),  
567 and now this is available to capture-recapture practitioners.

#### 568 4.1. Inference under boosting

569 One potential problem of boosted capture-recapture models is the new thinking required to understand  
570 what it is and how to describe its results. With origins in machine-learning, such algorithmic inference  
571 procedures may seem incomprehensible to ecologists: they may begrudge the lack of familiar inference tools  
572 like  $p$ -values and 95%CI (although, these are frequently misused; Hoekstra et al., 2014; Morey et al., 2016)  
573 or AIC weights. I offer two ways to understand boosting: comparison with other regularizers, and as a type  
574 of multi-model inference optimized for prediction.

575 In univariate analyses, boosting has some relationships to other procedures (see Meir & Rätsch, 2003, for  
576 an overview). For linear-models with Gaussian error, component-wise boosting is generally equivalent to the  
577 Lasso (Efron et al., 2004; Bühlmann & Hothorn, 2007). The Lasso can be viewed as simultaneously optimizing  
578 a goodness-of-fit term (i.e., a loss function) and a *penalty* on model complexity (the  $\ell_1$ -norm on regression  
579 coefficients). This form should be immediately familiar to most ecologists: the AIC also has a goodness-of-  
580 fit term and a fixed-penalty on model complexity ( $-2\ell_0$ -norm of regression coefficients). Hooten & Hobbs  
581 (2015) unify these ideas in a Bayesian framework: regularization is merely a strong prior disfavouring model  
582 complexity; more formally, regularized risk minimization is equivalent to Bayesian Maximum A-posteriori  
583 Probability (MAP) estimation (Murphy, 2012a), when the loss function is the negative log-likelihood. This is  
584 a helpful perspective, because inasmuch as capture-recapture practitioners are turning to Bayesian solutions  
585 under sparse data (Schofield et al., 2009; Schofield & Barker, 2011; Rankin et al., 2014, 2016), the CJSboost  
586 framework is allied and should be seriously considered. The above equivalences are more difficult to motivate  
587 using quixotic base-learners like CART-like trees, but which otherwise have great empirical performance  
588 under complex interactions and non-linear associations.

589 A second view of boosting is as an ensemble of many small models, like model-averaging. The terminology  
590 of a “learner” hails from its machine-learning origins, but base-learners are really just familiar analytic  
591 techniques commonly used for standalone modelling, like Ordinary Least Squares regression or CART. The  
592 influence of any one model is *weighted* according to the step-wise gradient descent procedure known as  
593 Boosting. Consider the case of Ordinary Least Squares base-learners: under extreme regularization ( $m=1$ ),  
594 the boosted estimates are the MLE of a simple intercept model (e.g.  $\hat{\phi}(\cdot)\hat{p}(\cdot)$ ). At weaker regularization  
595  $m \rightarrow \infty$ , the estimates tend to the MLEs of the fully-saturated model (Mayr et al., 2012). In between  
596 these extremes, at  $m_{\text{stop}} = m_{\text{CV}}$ , the estimates are *shrunk*, and somewhat qualitatively similar to AICc  
597 model-averaging. The size of the ensemble and its complexity is governed by predictive performance (through

598 cross-validation or bootstrap-validation). Thus, the resulting multimodel prediction function is that which  
599 minimizes the expected loss, and is therefore constrained from over-fitting. Unsurprisingly, the estimates  
600 have a slight downward bias but are more stable across outliers and different realizations of the data (i.e.  
601 favouring low-variance in the classic “bias-variance” trade-off; Bühlmann & Yu, 2003; Murphy, 2012a).

602 But what can one say about “significance” or “biological importance”? The answer is the interpretation  
603 of the additive coefficients (assuming they are similarly scaled): coefficients with the largest absolute values  
604 are the most influential on survival or capture probability. Using bootstrap stability-selection, we can also  
605 use approximate posterior inclusion probabilities as a type of uncertainty statistic: covariates/base-learners  
606 with high inclusion probabilities are probably more important; covariates with low inclusion probabilities  
607 ( $< 0.5$ ) are probably not that important. Probabilities lead to straight-forward inference. The stability  
608 paths (Figure 8) may also help visually discriminate between important covariates and noisy non-influential  
609 covariates, as suggested by Meinshausen & Bühlmann (2010): they notice a visual pattern whereby the  
610 true-model covariates enter the ensemble earlier and peel away from the unimportant covariates.

611 The above interpretations are hardly more difficult than understanding the AIC and model-averaging. In  
612 the applied ecological literature, there are few authors who formally justify a preference for the AIC versus  
613 other regularization and prediction techniques. Neither do ecologists seem to weigh in on philosophical  
614 arguments in favour of a prediction-optimal model versus a sparse model. Such matters are confused by a  
615 literature that is unclear about the underlying justification for AIC weighting and averaging (compare, for  
616 example, statements by Burnham, 2004, vs Raftery, 1995 and Hooten & Hobbs, 2015, about AIC weights  
617 as model probabilities). Commonly, ecologists cite “model parsimony” and Kullback-Leibler divergence as a  
618 justification for the AIC. This particular view of parsimony, however, favours certain outcomes.

619 Burnham (2004) offers a formal defence of the AIC and AIC model-averaging based on a notion of  
620 covariate “tapering”: the view that a response variable should theoretically have many small influences,  
621 possibly infinite, and our analyses should increasingly reveal more of these minor influences as we collect  
622 more data. They argue that natural phenomena are not “sparse”, unlike the systems studied by computer  
623 scientists, nor is there ever a “true model” (an oxymoron). This view is echoed by Vrieze (2012). The tapered  
624 worldview seems compelling for analyzing complex biological systems, where everything influences everything  
625 else. It is also, conveniently, the scenario in which the AIC and LOOCV are asymptotically prediction optimal  
626 and model-selection consistent (Shao, 1993, 1997; Burnham, 2004; Vrieze, 2012).

#### 627 *4.2. Tapering vs sparsity*

628 Nonetheless, I offer four arguments for capture-recapture methods to be more conservative. First, in  
629 an era of “Big Data” (geo-spatial, genetic, bio-logging, etc.) analysts increasingly have access to dozens of  
630 inventive potential covariates, many of which are different operationalizations of the same physical phenomena  
631 (e.g., consider the many ways one can measure Sea-Surface Temperature at multiple space-time scales).  
632 This Big Data deluge requires sparser discrimination among covariates, and if not, may encourage fishing  
633 for significance. Second, in an era of questionable scientific reproducibility (Simmons et al., 2011; Yaffe,

2015), we need better control on False Discoveries (among other things). This is a huge challenge, because from an optimal-prediction perspective, a False Rejection is much more costly to the expected loss than a shrunken False Discovery (Shao, 1993), thus making procedures overly liberal, including both the AIC and  $\ell_1$  regularizers. Third, there may be structural correlations in capture-recapture procedures that strongly favour certain outcomes, and which may preclude any hope for sparse, model-selection consistent estimates. I offer no theory to back this claim, but based on high-dimensional simulations, this study reveals high posterior inclusion probabilities for  $p(t)$  models (even when it is not the true model), as well as for covariates which are significant in one component, but not the other. This is likely not a feature of CJSboost, but a more widespread capture-recapture phenomenon (see Bailey et al., 2010 and Rankin et al., 2016, for problems of partial-identifiability of parameter estimates in capture-recapture). It can be expected to be more severe under low-detection probabilities. Fourth, in the author’s experience, the AIC/AICc seems to favour over-parametrized models that would be inadmissible under a Bayesian or a prediction paradigm, such as 100% survival and (the more ambiguous) 100% capture probability. Here, shrinkage on extreme values under regularization is similar to a Bayesian weak prior against boundary values.

To be clear, prediction-optimal  $\ell_1$  regularization, like L2boosting and the Lasso, are not very sparse, nor are they model-selection consistent (Meinshausen & Bühlmann, 2006; Zou, 2006; Bühlmann & Hothorn, 2010). They do, however, have more shrinkage on complexity than the AICc (Shao, 1997; Bühlmann & Hothorn, 2007) and AICc model-averaging, which is demonstrated in this study through simulation and an analysis of a real dataset. For more sparse model selection, the technique of bootstrapped stability selection (Meinshausen & Bühlmann, 2010; Murphy, 2012c) can be used to hard-threshold covariates which have low posterior inclusion probabilities ( $\lesssim 0.8-0.95$ ).

#### 4.3. Multimodel inference: build-up or post-hoc?

A boosted ensemble is built from the simplest intercept model and then “grows” more complex in a step-wise manner. This is the reverse of many multimodel inference techniques that do *post-hoc* weighting of models, such as AIC model-averaging and Bayesian model-averaging. However, the *post-hoc* approach becomes unmanageable with just a few covariates and parameters, given the combinatorial explosion in the number of plausible fixed-effect models. There is a risk that well-intentioned researchers will take short-cuts, such as a step-wise search strategy (Pérez-Jorge et al., 2016; Taylor et al., 2016), which may be susceptible to local-minima.

In conventional boosting, use of a convex loss function ensures that the gradient descent does not get stuck in a local minima. For non-convex problems, such as gamboostLSS (Mayr et al., 2012) and CJSboost, forced weakness/constraints on base-learners makes the problem more defined, but inevitably the start-values will dictate the direction of the gradient descent. However, for CJS and most capture-recaptures models, there is usually a well-defined intercept-only model that can serve as a principled way to initialize the predictions, such that if a unique MLE exists for the fully-saturated model, the boosting algorithm will reach it as  $m \rightarrow \infty$ . If there are parameter non-identifiability issues (such as for  $\{\phi_T, p_T\}$ ), early stopping will ensure that the

670 shrinkage is in the direction of the intercept-only model. Or, classic constraints can be imposed within the  
671 base-learners, such as fixing  $\phi_T = \phi_{T-1}$ .

#### 672 4.4. Extensions and future considerations

673 This study is merely the first step in developing and introducing boosting for HMM and capture-recapture.  
674 Many of the properties which hold for univariate Component-wise Boosting will need theoretical and empirical  
675 validation. Many questions arise, for example, how do the selection properties vary by sample-size, especially  
676 in reference to BIC and AIC model-averaging? How sensitive are the results to low detection probabilities?  
677 Does the EM technique and/or the MC technique generalize to multi-state models? How important is tuning  
678 both hyperparameters  $m$  and  $\nu$ ? Does the algorithm always reach the MLE of the fully-saturated model  
679 as  $m \rightarrow \infty$  and under what conditions does it fail? Is CJSboost and AICc-selection minimax optimal for  
680 mark-recapture?

681 By validating the boosting technique for a simple open-population model, this study paves the way for  
682 more popular capture-recapture models, such as POPAN and the PCRD, which have more model param-  
683 eters in the likelihood function, like temporary-migration processes. With more parameters, the boosting  
684 algorithms will require more efficient ways of tuning hyperparameters. See Appendix B.2 for ideas in this  
685 regard.

686 One major benefit of the CJSboost framework is its extensibility. It can easily accommodate phenomena  
687 such as individual heterogeneity, spatial capture-recapture and cyclic-splines. These are possible because  
688 the CJSboost code is written for compatibility with the `mboost` family of R packages, and leverages their  
689 impressive variety of base-learners (Bühlmann & Hothorn, 2007; Hofner et al., 2012). For example, the  
690 `brandom` base-learner can accommodate individual random effects for addressing individual heterogeneity in  
691 a manner similar to Bayesian Hierarchical models (Rankin et al., 2016). Kernels (`brad`) and spatial splines  
692 (`bspatial`) can be used for smooth spatial effects (Kneib et al., 2009; Hothorn et al., 2010; Tyne et al., 2015)  
693 offering an entirely new framework for spatial capture-recapture. The largest advantage is that users can  
694 add these extensions via the R formula interface, rather than having to modify deep-level code. CJSboost,  
695 therefore, offers a unified framework for many types of capture-recapture ideas that would otherwise require  
696 many different analytical paradigms to study the same suite of phenomena.

## 697 5. Conclusions

- 698 1. Boosting, the regularized gradient-descent and ensemble algorithm from machine learning, can be ap-  
699 plied to capture-recapture by reformulating the models as Hidden Markov Models, and interweaving an  
700 Expectation-Maximization E-step within each boosting iteration. An alternative boosting algorithm,  
701 based on stochastic imputation of HMM latent states, yields approximately equivalent estimates.
- 702 2. Boosting negotiates the “bias-variance” trade-off (for minimizing an expected loss) by incurring a slight  
703 bias in all coefficients, but yields estimates that are more stable to outliers and over-fitting, across



704 multiple realizations of the data. In contrast, Maximum Likelihood estimates are unbiased, but are  
705 highly variable.

706 3. CJSboost allows for powerful learners, such as recursive-partitioning trees (e.g., CART) for automatic  
707 variable-selection, interaction detection, and non-linearity. This flexibility seems to come at a cost of  
708 slightly more conservative estimates (if the underlying true model is linear).

709 4. Both AICc model-selection and boosting are motivated by good predictive performance: minimizing an  
710 expected loss, or generalization error. When using least-squares or CART-like base-learners, the esti-  
711 mates from CJSboost are qualitatively similar to AICc model-averaging, but with increased shrinkage  
712 on coefficients.

713 5. CJSboost seems to perform very well in high-dimensional model selection problems, with an ability to  
714 recover a small set of influential covariates. Typically, there is a small and non-zero weight on some  
715 unimportant covariates (especially  $p(t)$  base-learners). This pattern is consistent with the performance  
716 of univariate component-wise boosting and other  $\ell_1$  regularizers.

717 6. If the goal of a capture-recapture analysis is not prediction, but to recover a sparse “true model”, then  
718 CJSboosted models can be hard-thresholded via stability-selection. Hard-thresholded CJSboost models  
719 show some promise towards model-selection consistency and oracle-properties, but there may be some  
720 structural correlations in capture-recapture likelihoods that make this generally untrue.

## 721 6. Acknowledgements

722 I would like to thank Professor Sayan Mukherjee for inspiration during the Duke University course Prob-  
723 abilistic Machine Learning and giving this project an initial “thumbs up”.

## 724 7. Works Cited

725 Akaike, H. (1974). A new look at the statistical model identification. *Automatic Control, IEEE Transactions*  
726 *on*, 19, 716 – 723. doi:10.1109/TAC.1974.1100705.

727 Anderson, D. R., Burnham, K. P., & Thompson, W. L. (2000). Null hypothesis testing: problems, prevalence,  
728 and an alternative. *Journal of Wildlife Management*, 64, 912–923. doi:10.2307/3803199.

729 Bach, F. R. (2008). Bolasso: model consistent Lasso estimation through the bootstrap. In *Proceedings of the*  
730 *25th International Conference on Machine Learning ICML '08* (pp. 33–40). New York, NY, USA: ACM.  
731 doi:10.1145/1390156.1390161.

732 Bailey, L. L., Converse, S. J., & Kendall, W. L. (2010). Bias, precision, and parameter redundancy in complex  
733 multistate models with unobservable states. *Ecology*, 91, 1598–1604. doi:10.1890/09-1633.1.

734 Breiman, L. (1998). Arcing classifier (with discussion and a rejoinder by the author). *The Annals of Statistics*,  
735 26, 801–849. doi:10.1214/aos/1024691079.

- 736 Breiman, L. (1999). Prediction games and Arcing algorithms. *Neural Computation*, *11*, 1493–1517. doi:10.  
737 1162/089976699300016106.
- 738 Buckland, S. T., Burnham, K. P., & Augustin, N. H. (1997). Model selection: an integral part of inference.  
739 *Biometrics*, *53*, 603–618. doi:10.2307/2533961.
- 740 Burnham, K., & Anderson, D. (2014). P values are only an index to evidence: 20th-vs. 21st-century statistical  
741 science. *Ecology*, *95*, 627–630. doi:10.1890/13-1066.1.
- 742 Burnham, K. P. (2004). Multimodel inference: understanding AIC and BIC in model selection. *Sociological*  
743 *Methods & Research*, *33*, 261–304. doi:10.1177/0049124104268644.
- 744 Burnham, K. P., & Overton, W. S. (1978). Estimation of the size of a closed population when capture  
745 probabilities vary among animals. *Biometrika*, *65*, 625 – 633. doi:10.1093/biomet/65.3.625.
- 746 Bühlmann, P., & Hothorn, T. (2007). Boosting algorithms: regularization, prediction and model fitting.  
747 *Statistical Science*, *22*, 477–505. doi:10.1214/07-STS242.
- 748 Bühlmann, P., & Hothorn, T. (2010). Twin Boosting: improved feature selection and prediction. *Statistics*  
749 *and Computing*, *20*, 119–138. doi:10.1007/s11222-009-9148-5.
- 750 Bühlmann, P., & Yu, B. (2003). Boosting with the L2 loss: regression and classification. *Journal of the*  
751 *American Statistical Association*, *98*, 324–339. doi:10.1198/016214503000125.
- 752 Carothers, A. (1973). The effects of unequal catchability on Jolly-Seber estimates. *Biometrics*, *29*, 79–100.  
753 doi:10.2307/2529678.
- 754 Efron, B., Hastie, T., Johnstone, I., & Tibshirani, R. (2004). Least angle regression. *The Annals of Statistics*,  
755 *32*, 407–499. doi:10.1214/009053604000000067.
- 756 Elith, J., Leathwick, J. R., & Hastie, T. (2008). A working guide to boosted regression trees. *The Journal*  
757 *of Animal Ecology*, *77*, 802–813. doi:10.1111/j.1365-2656.2008.01390.x.
- 758 Friedman, J., Hastie, T., & Tibshirani, R. (2000). Additive logistic regression: a statistical view of boosting  
759 (with discussion). *The Annals of Statistics*, *28*, 337–374. doi:10.1214/aos/1016218223.
- 760 Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of Statistics*,  
761 *29*, 1189–1232.
- 762 Hand, D. J., & Vinciotti, V. (2003). Local versus global models for classification problems: fitting models  
763 where it matters. *The American Statistician*, *57*, 124–131. doi:10.1198/0003130031423.
- 764 Hoekstra, R., Morey, R. D., Rouder, J. N., & Wagenmakers, E.-J. (2014). Robust misinterpretation of  
765 confidence intervals. *Psychonomic Bulletin & Review*, *21*, 1157–1164. doi:10.3758/s13423-013-0572-3.

- 766 Hofner, B., Kneib, T., & Hothorn, T. (2014). A unified framework of constrained regression. *Statistics and*  
767 *Computing*, *26*, 1–14. doi:10.1007/s11222-014-9520-y.
- 768 Hofner, B., Mayr, A., Robinzonov, N., & Schmid, M. (2012). *Model-based Boosting in R: A Hands-on*  
769 *Tutorial Using the R Package mboost*. Technical Report 120 Department of Statistics, Ludwig-Maximilians-  
770 Universität Munich. URL: <http://epub.ub.uni-muenchen.de/12754/>.
- 771 Hooten, M., & Hobbs, N. T. (2015). A guide to Bayesian model selection for ecologists. *Ecological Monographs*,  
772 *85*, 3–28. doi:10.1890/14-0661.1.
- 773 Hothorn, T., Hornik, K., & Zeileis, A. (2006). Unbiased recursive partitioning: a conditional inference frame-  
774 work. *Journal of Computational and Graphical Statistics*, *15*, 651–674. doi:10.1198/106186006X133933.
- 775 Hothorn, T., Müller, J., Schröder, B., Kneib, T., & Brandl, R. (2010). Decomposing environmental, spatial,  
776 and spatiotemporal components of species distributions. *Ecological Monographs*, *81*, 329–347. doi:10.  
777 1890/10-0602.1.
- 778 Hunt, T., Bejder, L., Allen, S. J., Rankin, R. W., Hanf, D., & Parra, G. J. (2016). Demographic characteristics  
779 of Australian humpback dolphins reveal important habitat toward the south-western limit of their range.  
780 *Endangered Species Research*, in review.
- 781 Hutchinson, R., Liu, L., & Dietterich, T. (2011). Incorporating boosted regression trees into ecological latent  
782 variable models. In W. Burgard, & D. Roth (Eds.), *Proceedings of the Twenty-Fifth AAAI Conference on*  
783 *Artificial Intelligence* (pp. 1343–1348). Association for the Advancement of Artificial Intelligence.
- 784 Johnson, J. B., & Omland, K. S. (2004). Model selection in ecology and evolution. *Trends in Ecology &*  
785 *Evolution*, *19*, 101–108. doi:10.1016/j.tree.2003.10.013.
- 786 Kearns, M., & Valiant, L. (1994). Cryptographic limitations on learning Boolean formulae and finite au-  
787 tomata. *Journal of the ACM*, *41*, 67–95. doi:10.1145/174644.174647.
- 788 Kneib, T., Hothorn, T., & Tutz, G. (2009). Variable selection and model choice in geoaddivitive regression  
789 models. *Biometrics*, *65*, 626–634. doi:10.1111/j.1541-0420.2008.01112.x.
- 790 Laake, J. L. (2013). *RMark: An R Interface for Analysis of Capture-Recapture Data with MARK*. AFSC  
791 Processed Report 2013-01 Alaska Fisheries Science Center, NOAA, National Marine Fisheries Service  
792 Seattle, WA, USA.
- 793 Lebreton, J.-D., Burnham, K. P., Clobert, J., & Anderson, D. R. (1992). Modeling survival and testing  
794 biological hypotheses using marked animals: a unified approach with case studies. *Ecological Monographs*,  
795 *62*, 67–118. doi:10.2307/2937171.

- 796 Mayr, A., Fenske, N., Hofner, B., Kneib, T., & Schmid, M. (2012). Generalized additive models for location,  
797 scale and shape for high dimensional data—a flexible approach based on boosting. *Journal of the Royal*  
798 *Statistical Society: Series C (Applied Statistics)*, *61*, 403–427. doi:10.1111/j.1467-9876.2011.01033.x.
- 799 Meinshausen, N., & Bühlmann, P. (2006). High-dimensional graphs and variable selection with the Lasso.  
800 *The Annals of Statistics*, *34*, 1436–1462. doi:10.1214/009053606000000281.
- 801 Meinshausen, N., & Bühlmann, P. (2010). Stability selection. *Journal of the Royal Statistical Society: Series*  
802 *B (Statistical Methodology)*, *72*, 417–473. doi:10.1111/j.1467-9868.2010.00740.x.
- 803 Meir, R., & Rätsch, G. (2003). An introduction to boosting and leveraging. *Lecture Notes in Computer*  
804 *Science*, *2600*, 118–183. doi:10.1007/3-540-36434-X\_4.
- 805 Morey, R. D., Hoekstra, R., Rouder, J. N., Lee, M. D., & Wagenmakers, E.-J. (2016). The fallacy of  
806 placing confidence in confidence intervals. *Psychonomic Bulletin & Review*, *23*, 103–123. doi:10.3758/  
807 s13423-015-0947-8.
- 808 Mukherjee, S., Rifkin, R., & Poggio, T. (2003). Regression and Classification with Regularization. In D. D.  
809 Denison, M. H. Hansen, C. C. Holmes, B. Mallick, & B. Yu (Eds.), *Nonlinear estimation and classification*  
810 *Lecture Notes in Statistics* (pp. 111–128). New York, NY, USA: Springer.
- 811 Murphy, K. P. (2012a). Frequentist Statistics. In *Machine Learning: A Probabilistic Approach Adaptive*  
812 *computation and machine learning series* (pp. 191–216). Cambridge, MA, USA: MIT Press.
- 813 Murphy, K. P. (2012b). Markov and hidden Markov models. In *Machine Learning: A Probabilistic Approach*  
814 *Adaptive computation and machine learning series* (pp. 589–630). Cambridge, MA, USA: MIT Press.
- 815 Murphy, K. P. (2012c). Sparse Linear Models. In *Machine Learning: A Probabilistic Approach Adaptive*  
816 *computation and machine learning series* (pp. 421–478). Cambridge, MA, USA: MIT Press.
- 817 Oppel, S., Strobl, C., & Huettmann, F. (2009). *Alternative methods to quantify variable importance in*  
818 *ecology*. Technical Report 65 Department of Statistics, Ludwig-Maximilians-Universität Munich.
- 819 Pérez-Jorge, S., Gomes, I., Hayes, K., Corti, G., Louzao, M., Genovart, M., & Oro, D. (2016). Effects  
820 of nature-based tourism and environmental drivers on the demography of a small dolphin population.  
821 *Biological Conservation*, *197*, 200–208. doi:10.1016/j.biocon.2016.03.006.
- 822 Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition.  
823 *Proceedings of the IEEE*, *77*, 257–286. doi:10.1109/5.18626.
- 824 Raftery, A. E. (1995). Bayesian model selection in social research. *Sociological methodology*, *25*, 111–164.
- 825 Rankin, R. W., Maldini, D., & Kaufman, G. D. (2014). Bayesian estimate of Australian humpback whale calv-  
826 ing intervals under sparse resighting rates: 1985 – 2009. *Journal of Cetacean Research and Management*,  
827 *13*, 109–121.

- 828 Rankin, R. W., Nicholson, K. E., Allen, S. J., Krützen, M., Bejder, L., & Pollock, K. H. (2016). A full-capture  
829 Hierarchical Bayesian model of Pollock’s Closed Robust Design and application to dolphins. *Frontiers in*  
830 *Marine Science*, 3. doi:10.3389/fmars.2016.00025.
- 831 Robinzonov, N. (2013). *Advances in boosting of temporal and spatial models*. Doctoral Thesis LMU München:  
832 Fakultät für Mathematik, Informatik und Statistik Munich. URL: [http://nbn-resolving.de/urn:nbn:](http://nbn-resolving.de/urn:nbn:de:bvb:19-153382)  
833 [de:bvb:19-153382](http://nbn-resolving.de/urn:nbn:de:bvb:19-153382).
- 834 Schapire, R. E. (1990). The strength of weak learnability. *Machine Learning*, 5, 197–227. doi:10.1023/A:  
835 1022648800760.
- 836 Schmid, M., & Hothorn, T. (2008a). Boosting additive models using component-wise P-Splines. *Computa-*  
837 *tional Statistics & Data Analysis*, 53, 298–311. doi:10.1016/j.csda.2008.09.009.
- 838 Schmid, M., & Hothorn, T. (2008b). Flexible boosting of accelerated failure time models. *BMC Bioinform-*  
839 *atics*, 9, 269. doi:10.1186/1471-2105-9-269.
- 840 Schmid, M., Potapov, S., Pfahlberg, A., & Hothorn, T. (2010). Estimation and regularization techniques  
841 for regression models with multidimensional prediction functions. *Statistics and Computing*, 20, 139–150.  
842 doi:10.1007/s11222-009-9162-7.
- 843 Schmid, M., Wickler, F., Maloney, K. O., Mitchell, R., Fenske, N., & Mayr, A. (2013). Boosted Beta  
844 Regression. *PLoS ONE*, 8. doi:10.1371/journal.pone.0061623.
- 845 Schofield, M. R., & Barker, R. J. (2011). Full open population capture–recapture models with individual  
846 covariates. *Journal of Agricultural, Biological, and Environmental Statistics*, 16, 253–268. doi:10.1007/  
847 s13253-010-0052-4.
- 848 Schofield, M. R., Barker, R. J., & Mackenzie, D. I. (2009). Flexible hierarchical mark-recapture modeling for  
849 open populations using WinBUGS. *Environmental and Ecological Statistics*, 16, 369–387. doi:10.1007/  
850 s10651-007-0069-1.
- 851 Shao, J. (1993). Linear Model selection by cross-validation. *Journal of the American Statistical Association*,  
852 88, 486–494. doi:10.1080/01621459.1993.10476299.
- 853 Shao, J. (1997). An asymptotic theory for linear model selection. *Statistica Sinica*, 7, 221–242.
- 854 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology undisclosed flexibility in  
855 data collection and analysis allows presenting anything as significant. *Psychological Science*, 22, 1359–1366.  
856 doi:10.1177/0956797611417632.
- 857 Stone, M. (1977). An asymptotic equivalence of choice of model by cross-validation and Akaike’s criterion.  
858 *Journal of the Royal Statistical Society. Series B (Methodological)*, 39, 44–47.

- 859 Taylor, A. R., Schacke, J. H., Speakman, T. R., Castleberry, S. B., & Chandler, R. B. (2016). Factors related  
860 to common bottlenose dolphin (*Tursiops truncatus*) seasonal migration along South Carolina and Georgia  
861 coasts, USA. *Animal Migration*, *3*. doi:10.1515/ami-2016-0002.
- 862 Tibshirani, R. (2011). Regression shrinkage and selection via the lasso: a retrospective: Regression Shrinkage  
863 and Selection via the Lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,  
864 *73*, 273–282. doi:10.1111/j.1467-9868.2011.00771.x.
- 865 Tyne, J. A., Johnston, D. W., Rankin, R., Loneragan, N. R., & Bejder, L. (2015). The importance of spinner  
866 dolphin (*Stenella longirostris*) resting habitat: implications for management. *Journal of Applied Ecology*,  
867 *52*, 621–630. doi:10.1111/1365-2664.12434.
- 868 Vrieze, S. I. (2012). Model selection and psychological theory: A discussion of the differences between the  
869 Akaike information criterion (AIC) and the Bayesian information criterion (BIC). *Psychological Methods*,  
870 *17*, 228–243. doi:10.1037/a0027127.
- 871 Ward, G., Hastie, T., Barry, S., Elith, J., & Leathwick, J. R. (2009). Presence-only data and the EM  
872 algorithm. *Biometrics*, *65*, 554–563. doi:10.1111/j.1541-0420.2008.01116.x.
- 873 White, G. C., & Burnham, K. P. (1999). Program MARK: survival estimation from populations of marked  
874 animals. *Bird Study*, *46*, S120–S139. doi:10.1080/00063659909477239.
- 875 Yaffe, M. B. (2015). Reproducibility in science. *Science Signaling*, *8*, EG5. doi:10.1126/scisignal.aaa5764.
- 876 Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*,  
877 *101*, 1418–1429. doi:10.1198/016214506000000735.

## 878 APPENDICES

### 879 Appendix A. Algorithms for Filtering and Sampling HMM Latent States

880 The CJSboost algorithm depends on conditional independence of data pairs  $(y_{i,t}, X_{i,t})$  for individuals  $i$   
881 in capture period  $t$ , in order to estimate the negative-gradient in the descent algorithm. This is possible if  
882 we impute information about the latent state sequences  $z$  for pairs of capture periods at  $t$  and  $t-1$ . The  
883 two CJSboost algorithms, CJSboost-EM and CJSboost-MC, achieve this same idea with two different, but  
884 related, techniques. In both cases, we will use a classic “forwards-backwards” messaging algorithm to gain  
885 information about the probability distribution of the latent state sequences. In CJSboost-EM, we calculate  
886 the *two-slice marginal probabilities*  $p(z_{t-1} = u, z_t = v | \mathbf{y}_{1:T}, \phi, p)$ , per boosting iteration; in CJSboost-MC, we  
887 will *sample*  $\mathbf{z}$  from its posterior distribution  $\pi(\mathbf{z}_{1:T} | \mathbf{y}_{1:T}, \phi, p)$ . See Rabiner (1989) and Murphy (2012b) for  
888 accessible tutorials.

889 Both algorithms use a forwards algorithm and backwards algorithm. We will drop the indices  $i$ , and focus  
890 on the capture history of a single individual.  $\mathbf{y}$  is the time-series of binary outcomes of length  $T$ .  $\mathbf{z}$  is a

891 vector of latent states  $z \in \{\text{dead}, \text{alive}\}$ . We condition on an individual's first capture at time  $t = t^0$ , and are  
 892 only concerned with the sequence  $\mathbf{z}_{t^0:T}$ . Survival from step  $t-1$  to  $t$  is  $\phi_t$ . Conditional on  $z_t$ , the capture  
 893 probabilities are  $p(y_t = 1|\text{alive}) = p_t$ , and  $p(y_t = 1|\text{dead}) = 0$ . In HMM notation, the CJS processes can be  
 894 presented as the following column-stochastic matrices:

$$\Phi_t = \begin{array}{c} \text{dead} \\ \text{alive} \end{array} \begin{array}{cc} \text{dead} & \text{alive} \\ \left( \begin{array}{cc} 1 & 1-\phi_t \\ 0 & \phi_t \end{array} \right) \end{array} \quad \Psi_t = \begin{array}{c} \text{no capture} \\ \text{capture} \end{array} \begin{array}{cc} \text{dead} & \text{alive} \\ \left( \begin{array}{cc} 1 & 1-p_t \\ 0 & p_t \end{array} \right) \end{array} \quad (\text{A.1})$$

895 In HMM parlance,  $\Phi$  is the Markovian transition process; we denote the probability  $p(z_t = u|z_{t-1} = u)$  as  
 896  $\Phi_t(u, v)$ .  $\Psi$  is the emission process governing capture probabilities; we denote the probability  $p(y_t = 1|z_t = v)$   
 897 as  $\Psi_t(v)$ .

#### 898 *Appendix A.1. Forward-algorithm*

899 The forward messaging algorithm involves the recursive calculation of  $\alpha_t(v)$ , per time  $t$  and state  $z_t = v$ .  
 900  $\alpha_t$  is the *filtered belief state* of  $z_t$  given all the observed information in  $\mathbf{y}$  from first capture  $t^0$  until  $t$ . Notice,  
 901 that for clarity, we drop the notation for conditioning on  $\phi$  and  $p$ , but these are always implied.

$$\begin{aligned} a_t(v) &:= p(z_t = v|\mathbf{y}_{t^0:t}) \\ &= \frac{1}{Z_t} p(y_t|z_t = v)p(z_t = v|\mathbf{y}_{t^0:t-1}) \\ &= \frac{1}{Z_t} p(y_t|z_t = v) \sum_u p(z_t = v|z_{t-1} = u)p(z_{t-1} = u|\mathbf{y}_{t^0:t-1}) \\ &= \frac{1}{Z_t} \Psi_t(v) \sum_u \Phi(u, v)\alpha_{t-1}(u) \\ Z_t &= \sum_v \left( \Psi_t(v) \sum_u \Phi(u, v)\alpha_{t-1}(u) \right), \quad \sum_v \alpha_t(v) = 1 \end{aligned} \quad (\text{A.2})$$

902 The algorithm is initialized at time  $t^0$  (an individual's first capture) with  $\alpha_{t^0}(\text{alive}) = 1$ . Conditional on the  
 903 values of  $\alpha_t(v)$  for all  $v$ , one can proceed to calculate the next values of  $\alpha_{t+1}(v)$ , and so on, until  $t = T$ .

#### 904 *Appendix A.2. Backwards-algorithm*

905 Messages are passed backwards in a recursive algorithm starting at  $t = T$  and moving backwards until  
 906  $t = t^0$ , the first-capture period, while updating entries in  $\beta_t(v)$ .

$$\begin{aligned} \beta_{t-1}(u) &:= p(\mathbf{y}_{t:T}|z_{t-1} = u) \\ &= \sum_v p(\mathbf{y}_{t+1:T}|z_t = v)p(y_t|z_t = v)p(z_t = v|z_{t-1} = u) \\ &= \sum_v \beta_t(v)\Psi_t(v)\Phi_t(u, v) \end{aligned} \quad (\text{A.3})$$

907 The algorithm is initialized  $\beta_T(\cdot) = 1$  for all states  $v$  (notice that the entries do not need to sum to 1).  
 908 Having calculated the backwards and forwards messages, we can now proceed to characterize the latent state  
 909 distributions.



910 *Appendix A.3. Two-slice marginal probabilities for Expectation-Maximization*

911 Expectation-Maximization is an iterative technique for maximizing a difficult objective function by work-  
 912 ing with an easy “complete-data” objective function  $\log p(y, z|\theta)$ . EM works by cycling through an M-step  
 913 and an E-step. In boosting-EM, the M-step corresponds to the usual update of the prediction vectors  
 914  $F_\theta^{(m)} = F_\theta^{(m-1)} + \nu \theta \hat{f}$  (conditional on  $z$ ), and are used to estimate  $\hat{\theta}$ . The E-step imputes expectations of the  
 915 latent states  $z$ , conditional on the data and current estimates of  $\hat{\theta}^{(m)}$ .

916 In the CJSboost-EM algorithm, we require expectations for the joint states  $(z_{t-1}, z_t)$ . We substitute in  
 917 the two-slice marginal probabilities  $p(z_{t-1}, z_t | \mathbf{y}_{t^0:T}, \phi, p)$ . These can be easily evaluated for a capture history  
 918  $\mathbf{y}_i$  using the outputs  $(\alpha, \beta)$  from the forward-backwards algorithm.

$$\begin{aligned} w_t(u, v) &:= p(z_{t-1} = u, z_t = v | \mathbf{y}_{t^0:T}) \\ &= \frac{1}{\xi_t} p(z_{t-1} | \mathbf{y}_{t^0:t-1}) p(z_t | z_{t-1}, \mathbf{y}_{t:T}) \\ &= \frac{1}{\xi_t} p(z_{t-1} | \mathbf{y}_{t^0:t-1}) p(y_t | z_t) p(\mathbf{y}_{t+1:T} | z_t) p(z_t | z_{t-1}) \\ &= \frac{1}{\xi_t} \alpha_{t-1}(u) \Psi_t(v) \beta_t(v) \Phi_t(u, v) \end{aligned} \tag{A.4}$$

$$\xi_t = \sum_u \sum_v \alpha_{t-1}(u) \Psi_t(v) \beta_t(v) \Phi_t(u, v), \quad \sum_u \sum_v w_t(u, v) = 1$$

919 The E-step is completed after evaluating the set  $\{w_{i,t}(\text{alive}, \text{alive}), w_{i,t}(\text{alive}, \text{dead}), w_{i,t}(\text{dead}, \text{dead})\}$ , for  
 920 each capture period  $t > t_i^0$  and for each individual capture history  $\{\mathbf{y}_i\}_{i=1}^n$ . This is an expensive operation;  
 921 computational time can be saved by re-evaluating the expectations every second or third boosting iteration  
 922  $m$ , which, for large  $m_{\text{stop}} > 100$  and small  $\nu$ , will have a negligible approximation error.

923 *Appendix A.4. Sampling state-sequences from their posterior*

924 For the CJSboost Monte-Carlo algorithm, we sample a latent state sequence  $\mathbf{z}_i$  from the posterior  
 925  $\pi(\mathbf{z}_{1:T} | \mathbf{y}_{1:T}, \phi, p)$ , for each individual  $i$  per boosting step. Conditional on the latent states, the negative-  
 926 gradients are easily evaluated and we can proceed to boost the estimates and descend the risk gradient.  
 927 However, because the algorithm is stochastic, we must avoid getting trapped in a local minima by sampling  
 928 many sequences (e.g.,  $S \approx 10 - 20$ ), thereby approximating the full posterior distribution of  $\mathbf{z}$ . Over all  $S$   
 929 samples, the average gradient will *probably* be in the direction of the global minima. For large  $m$  and small  
 930  $\nu$ , the approximation error is small.

The algorithm uses backwards-sampling of the posterior under the chain rule:

$$p(\mathbf{z}_{t^0:T} | \mathbf{y}_{t^0:T}) = p(z_T | \mathbf{y}_{t^0:T}) \prod_{t=T-1}^{t^0} p(z_t | z_{t+1}, \mathbf{y}_{t^0:T}) \tag{A.5}$$

931 We start with a draw at time  $t = T$ ,  $z_T^{(s)} \sim p(z_T = v | \mathbf{y}_{t^0:T}) = \alpha_T(v)$ , and condition earlier states on

932 knowing the next-step-ahead state, proceeding backwards until  $t = t^0$ .

$$\begin{aligned}
 z_t^{(s)} &\sim p(z_t = u | z_{t+1} = v, \mathbf{y}_{t^0:t}) \\
 &= \frac{p(z_t, z_{t+1} | \mathbf{y}_{t^0:t+1})}{p(z_{t+1} | \mathbf{y}_{t^0:t+1})} \\
 &\propto \frac{p(y_{t+1} | z_{t+1}) p(z_t, z_{t+1} | \mathbf{y}_{t^0:t})}{p(z_{t+1} | \mathbf{y}_{t^0:t+1})} \\
 &= \frac{p(y_{t+1} | z_{t+1}) p(z_{t+1} | z_t) p(z_t | \mathbf{y}_{t^0:t})}{p(z_{t+1} | \mathbf{y}_{t^0:t+1})} \\
 &= \frac{\Psi_{t+1}(v) \Phi_{t+1}(u, v) \alpha_t(u)}{\alpha_{t+1}(v)}
 \end{aligned} \tag{A.6}$$

933 Thus, knowing  $\alpha$ ,  $\beta$ ,  $\Phi$  and  $\Psi$ , we can easily generate random samples of  $\mathbf{z}$  that are drawn from its  
 934 posterior distribution. The backwards sampling step is repeated for each  $t > t_i^0$  capture period, for each  $s$   
 935 sequence, for each  $i$  capture history, for each  $m$  boosting iteration.

## 936 Appendix B. Tuning Hyperparameters $m$ and $\nu$

937 This section will present a simple work-flow for finding approximately optimal values of  $m_{stop}$ ,  $\nu_\phi$  and  
 938  $\nu_p$ . Our objective is to minimize the expected loss  $\mathcal{L}$ , or generalization error. We estimate  $\mathcal{L}$  through  $B$ -  
 939 times bootstrap-validation. For each  $b$  bootstrap, we create a CJSboost prediction function,  $G^{(b)}(X; m, \nu_\phi, \nu_p)$   
 940 which is trained on the bootstrapped data and is a function of the hyperparameters  $\nu_\phi$ ,  $\nu_p$  and  $m$ . We calculate  
 941 the holdout-out risk using the out-of-bootstrap  $b^c$  capture-histories and covariate data,  $(\mathbf{Y}^{(b^c)}, \mathbf{X}^{(b^c)})$ . The  
 942 average hold-out risk over  $B$  bootstraps,  $L_{cv}$ , is our objective to minimize.

$$\mathcal{L} \approx L_{cv} = \operatorname{argmin}_{m, \nu_\phi, \nu_p} \frac{1}{B} \sum_{b=1}^B L(\mathbf{Y}^{(b^c)}, G^{(b)}(\mathbf{X}^{(b^c)}; m, \nu_\phi, \nu_p))$$

943 For a given  $\nu_\phi$  and  $\nu_p$ , the hold-out risk can be monitored internally to the boosting algorithm for each  
 944 step  $m$ . Therefore, a single B-bootstrap run is all that is necessary to find the optimal  $m$ , given  $\nu_\phi$  and  
 945  $\nu_p$ . But since  $\nu_\phi$  and  $\nu_p$  are continuous, one must discretize the range of possible values and re-run separate  
 946 B-bootstrap-validation exercises per combination of  $\nu_\phi$  and  $\nu_p$ . This is very expensive, and one must accept  
 947 some approximation error.

### 948 Appendix B.1. Algorithm 1 for tuning $\nu$

949 For just two parameters, the pertinent quantity to optimize is the ratio  $\lambda = \frac{\nu_p}{\nu_\phi}$ , for a fixed mean  $\nu_\mu =$   
 950  $\frac{1}{2}(\nu_\phi + \nu_p)$ . Therefore, a univariate discrete set of  $\Lambda = \{\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(J)}\}$  can be searched for the smallest  
 951  $L_{cv}(\lambda)$ .

952 This is less daunting than it may seem, because the range of  $\lambda$  is practically bounded. For example,  
 953 if  $m_{stop} = 1000$  and  $\lambda = 100$ ,  $\phi$  is effectively shrunk to its intercept starting value, and higher values of  
 954  $\lambda$  have little effect. Also, Bühlmann & Yu (2003) suggest that the generalization error has a very shallow  
 955 minima around the optimal values of  $m$ , which means that our hyperparameters need only get within the

vicinity of their optimal values, rather than strict numerical convergence. Finally,  $L_{cv}(\lambda)$  is typically convex for varying  $\lambda$  (so long as the same bootstrap-weights are recycled for all new estimates of  $L_{cv}(\lambda)$ ). Therefore, we can employ any convex optimization algorithm for non-differentiable functions to iteratively search for the optimal  $\lambda$ . The thrust of any such algorithm is a multiplicative “stepping-out” procedure to quickly find the correct order of magnitude for  $\lambda$ . For example, starting a  $\lambda^{(0)} = 1$ , we need only 7 doubling steps to grow  $\lambda$  to  $128 \times \lambda^{(0)}$ ; further refinements will have little practical impact on the final model estimates.

An example algorithm is the following.

1. set  $\nu_\mu = 0.01$  and  $\lambda^{(0)} = 1$ ; generate the  $B$  bootstrap samples; initialize the set  $\Lambda = \{\lambda^{(0)}, \frac{1}{2}\lambda^{(0)}\}$ ;
2. for each  $\lambda$  in  $\Lambda$ , estimate  $L_{cv}(\lambda)$  and store the values in the list  $\mathbf{L} = \{L^{(0)}, \dots\}$ ;
3. for  $j$  in  $1:J$ , do:
  - (a) get the current best value for the ratio  $\lambda_{\min} = \underset{\lambda \in \Lambda}{\operatorname{argmin}} L_{cv}(\lambda)$
  - (b) estimate a new candidate  $\lambda^*$ :
    - if  $\lambda_{\min} = \min(\Lambda)$ , then  $\lambda^* = \frac{1}{2}\min(\Lambda)$ ;
    - else if  $\lambda_{\min} = \max(\Lambda)$ , then  $\lambda^* = 2 \cdot \max(\Lambda)$ ;
    - else  $\lambda^* = \lambda_{\min} + k \cdot \alpha$ , where  $k$  is the step direction and  $\alpha$  is the step size.
  - (c) calculate the shrinkage weights:  $\nu_\phi^{(j)} = \frac{2 \cdot \nu_\mu}{\lambda^* + 1}$ ;  $\nu_p^{(j)} = \lambda^* \cdot \nu_\phi^{(j)}$ ;
  - (d) perform bootstrap-validation to estimate  $L_{cv}^{(j)}(\lambda^*)$ ;
  - (e) append  $\Lambda \leftarrow \lambda^*$  and append  $\mathbf{L} \leftarrow L_{cv}^{(j)}$ ;

The algorithm continues until a pre-defined convergence criteria is met, or, practically, a maximum number of iterations is reached. The final values of  $\nu_\phi$ ,  $\nu_p$ , and  $m_{cv}$  are those which correspond to the minimum  $L_{cv} \in \mathbf{L}$ .

There are various convex optimization algorithms that differ in how to calculate the  $k$  and  $\alpha$ . In CJSboost, most of the optimization benefits occur during the “stepping-out” procedure, and so exact values of  $k$  and  $\alpha$  are less important, so long as they guarantee convergence. I suggest the following sub-algorithm (nested within step 3b above), which converges slowly but quickly rules out large chunks of bad values of  $\lambda$ .

1. Define the triplet set  $\Gamma$  composed of the current best estimate of  $\lambda_{\min}$  as well as the values just to the left and right, such that  $\lambda_{\min}^{-1} < \lambda_{\min} < \lambda_{\min}^{+1}$ ;
2. Sort the entries of  $\Gamma$  according to the order  $L_{cv}(\gamma^{(1)}) < L_{cv}(\gamma^{(2)}) < L_{cv}(\gamma^{(3)})$ ;
3. Estimate the step size and direction:
  - if  $\|\gamma^{(1)} - \gamma^{(2)}\| \geq \|\gamma^{(1)} - \gamma^{(3)}\|$ :
  - then  $\alpha = \frac{1}{2}\|\gamma^{(1)} - \gamma^{(2)}\|$  and  $k = \operatorname{sign}(\gamma^{(1)} - \gamma^{(2)})$ ;
  - else  $\alpha = \frac{1}{2}\|\gamma^{(1)} - \gamma^{(3)}\|$  and  $k = \operatorname{sign}(\gamma^{(1)} - \gamma^{(3)})$ ;
4.  $\lambda^* = \lambda_{\min} + k \cdot \alpha$

Typically seven or ten iterations are necessary in order to find suitable values of  $\lambda$ ,  $\nu_\phi$  and  $\nu_p$ . Unfortunately, this strategy is only for a two-parameter likelihood with a single ratio to optimize. For other capture-recapture models with more parameters (e.g., POPAN, PCRD), a different tuning strategy will be necessary.

992 *Appendix B.2. Algorithm 2 for tuning  $\nu$*

993 With more parameters in the capture-recapture likelihood, the number of necessary steps in algorithm  
994 1 will increase exponentially. I suggest a second iterative algorithm whose number of iterations may only  
995 increase linearly with the number of parameters. The principle of this second algorithm is based on the  
996 observation that when the ratio  $\frac{\nu_p}{\nu_\phi}$  is poorly optimized, then additional boosting steps along the gradient  
997  $\frac{\partial \ell}{\partial F_\theta}$  will result in *increases* in the holdout-risk, and will do so asymmetrically for  $F_\phi$  vs  $F_p$ . When  $\frac{\nu_p}{\nu_\phi}$  is  
998 optimized, the number of boosting steps which increase the hold-out risk will be roughly the same for  $p$  and  
999  $\phi$ , averaged over all bootstrap hold-out samples. I suggest using this ratio as an estimate of  $\hat{\lambda} = \frac{\nu_p}{\nu_\phi}$ .

1000 Call  $\Delta_\theta^{(m)}$  a boosting step along the partial derivative of  $\frac{\partial \ell}{\partial F_\theta}$  which successfully reduces the holdout-risk.

$$\hat{\lambda}^{(j)} = \hat{\lambda}^{(j-1)} Q \left( \frac{\sum_{m=1}^{m_k} \Delta_p^{(m)}}{\sum_{m=1}^{m_k} \Delta_\phi^{(m)}} \right) \quad (\text{B.1})$$

1001 where  $Q$  is a robust measure of central tendency over all  $B$  bootstraps (median, trimmed-mean), and  $m_k$   
1002 is some boosting step  $m_k > m_{cv}$ . The first estimate  $\hat{\lambda}^{(1)}$  is typically an underestimate, so the algorithm is  
1003 iterated, each time using the previous  $\hat{\lambda}^{(j-1)}$  for a current estimate of  $\nu_p$  and  $\nu_\phi$  with which to perform a  
1004 bootstrap-validation exercise, and then updating  $\hat{\lambda}^{(j)}$  by (B.1).  $\hat{\lambda}^{(J)}$  typically converges to a single value  
1005 within approximately 10 iterations.  $\hat{\lambda}^{(J)}$  is *not* the optimal  $\lambda$  as estimated by algorithm 1, but it is in the  
1006 vicinity (Figure B.9).

1007 Clearly, for just two parameters and one ratio, this second algorithm is not competitive with algorithm  
1008 1. But, when there are more than two parameters in the likelihood, this algorithm can simultaneously  
1009 estimate all pertinent ratios. Further refinements will be necessary, but simulations demonstrate that there  
1010 is information in the risk gradient trajectories that can help optimize the hyperparameters.

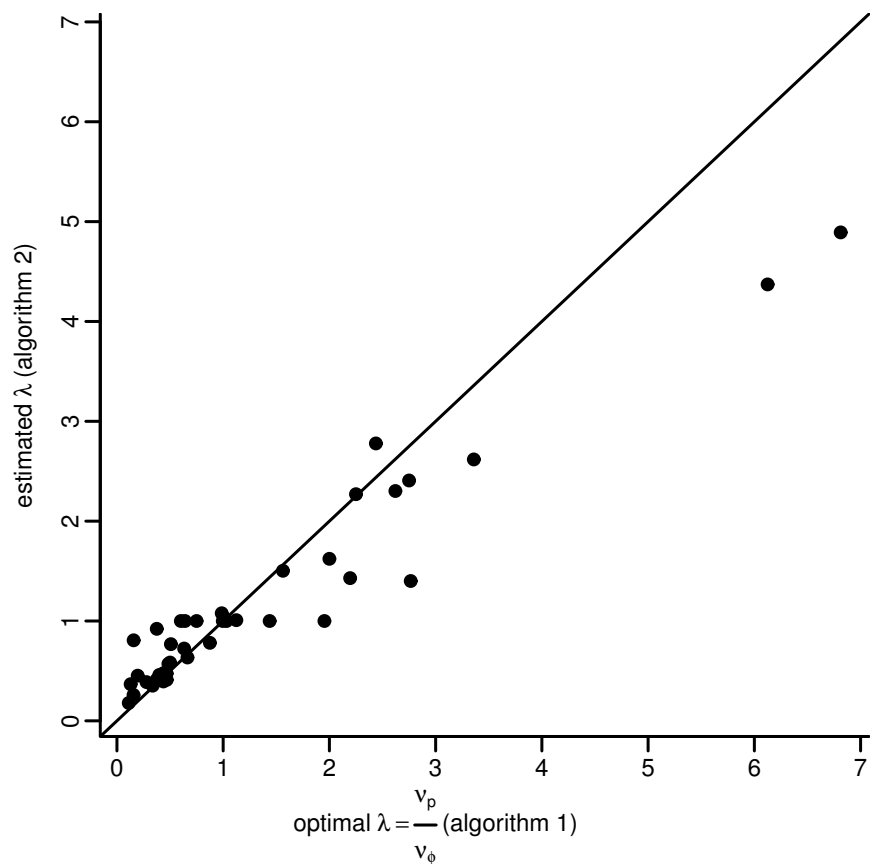


Figure B.9: Two algorithms for tuning the shrinkage weight hyperparameters  $\nu_\phi$  and  $\nu_p$ , and their ratio  $\lambda$ , in order to minimize the expected loss (estimated via bootstrap-validation). Forty simulations compare the two algorithms, where algorithm 1 is considered optimal.