

<sup>1</sup> RH: Model Selection for North American BBS

<sup>2</sup> **Model selection for the North American Breeding Bird  
3 Survey, with observations on BPIC and WAIC model  
4 selection criteria**

<sup>5</sup> WILLIAM A. LINK<sup>1,2</sup> , JOHN R. SAUER<sup>1</sup>, AND DANIEL K. NIVEN<sup>1</sup>

<sup>6</sup> *<sup>1</sup>USGS Patuxent Wildlife Research Center, Laurel, Maryland 20708, USA*

---

<sup>2</sup>E-mail: [wlink@usgs.gov](mailto:wlink@usgs.gov)

7        *Abstract.* The North American Breeding Bird Survey (BBS) provides data that can be  
8        used in complex, multi-scale analyses of population change, while controlling for scale-specific  
9        nuisance factors. Many alternative models can be fit to the data, but most model selection  
10        procedures are not readily applied to hierarchical models. Leave-one-out cross validation  
11        (LOOCV), in which relative model fit is assessed by omitting an observation and assessing  
12        the prediction of a model fit using the remainder of the data, provides a reasonable approach  
13        for assessing models, but is time consuming and not feasible to apply for all observations in  
14        large data sets. We report the first large-scale formal model selection for BBS data, applying  
15        LOOCV to stratified random samples of observations from BBS data. Our results are for 548  
16        species of North American birds, comparing the fit of 4 alternative models that differ in year-  
17        effect structures and in descriptions of extra-Poisson overdispersion. We use a hierarchical  
18        model among species to evaluate posterior probabilities that models are best for individual  
19        species. Models in which differences in year effects are conditionally independent (D models)  
20        were generally favored over models in which year effects are modeled by a slope parameter  
21        and a random year effect (S models), and models in which extra-Poisson overdispersion  
22        effects are independent and t-distributed (H models) tended to be favored over models  
23        where overdispersion was independent and normally distributed. Our conclusions lead us to  
24        recommend a change from the conventional S model to D and H models for the vast majority  
25        of species (544/548). Comparison of estimated population trends based on the favored model  
26        relative to the S model currently used for BBS summaries indicates no consistent differences  
27        in estimated trends. Of the 18 species that showed large differences in estimated trends  
28        between models, estimated trends from the default S model were more extreme, reflecting  
29        the influence of the slope parameter in that model for species that are undergoing large  
30        population changes. WAIC, a computationally simpler alternative to LOOCV, does not

<sup>31</sup> appear to be a reliable alternative to LOOCV.

<sup>32</sup> *Key words:* Bayesian analysis; hierarchical models; leave-one-out cross validation; model  
<sup>33</sup> selection; North American Breeding Bird Survey; WAIC

<sup>34</sup> **Introduction**

<sup>35</sup> The North American Breeding Bird Survey (BBS) provides count data used to monitor  
<sup>36</sup> population change for over 500 species of birds. The survey, in which data are collected  
<sup>37</sup> annually by skilled observers on preestablished survey routes, began in 1966 in the eastern  
<sup>38</sup> United States and has expanded in size and geographic extent (Sauer et al. 2013). By 2018,  
<sup>39</sup> the BBS database contained information from 5745 routes across the contiguous United  
<sup>40</sup> States, Alaska, and Canada, of which  $\sim$  3000 are surveyed annually. Additional survey routes  
<sup>41</sup> have been established in Mexico, but are not routinely analyzed due to limited data.

<sup>42</sup> The BBS is the only source of information on population change for most North American  
<sup>43</sup> bird species, and is the primary data source for conservation status assessments (Rosenberg et  
<sup>44</sup> al. 2016) and State of the Birds Reports (North American Bird Conservation Initiative 2016).

<sup>45</sup> The United States Geological Survey (USGS) publishes estimates of population change at  
<sup>46</sup> a variety of spatial scales, with results published via yearly website updates (e.g., Sauer et  
<sup>47</sup> al. 2014). The wide use of these results is evidenced by over 1300 published citations of the  
<sup>48</sup> USGS website in the last 20 years.

<sup>49</sup> BBS data present a challenge for analysts. The survey spans the continent, varies greatly in  
<sup>50</sup> consistency of coverage over space and time, and is conducted by thousands of observers that  
<sup>51</sup> vary in skills (Sauer et al. 2013). Change in observers is an important feature of the data:  
<sup>52</sup> the typical BBS observer provides 4 counts over 5 years of service on a BBS route. Observers'

53 participation duration on a route have 25th, 50th, 75th, and 90th percentiles of 1, 5, 12 and  
54 23 years (Link and Sauer 2016a). Bird species differ greatly in life history attributes and  
55 geographic distributions, and these differences also have consequences for statistical analysis.

56 To accommodate these complexities in the analysis, many analysts of BBS data use overdis-  
57 persed Poisson regression models, analyzed as Bayesian hierarchical models (e.g., Sauer et al.  
58 2011). Expected counts are functions of observer effects and spatially stratified year effects,  
59 which we now describe.

60 Observer effects are of two sorts, “among-observer” and “within-observer.” Among-observer  
61 effects describe variation in count rates for different observers under identical circumstances.  
62 Temporal changes in counts attributable to changes in the pool of observers are well-  
63 documented (Link and Sauer 1998); overlooking these effects biases estimates of population  
64 change, typically leading to positive biases in trends. Within-observer effects reflect temporal  
65 changes in an individual’s counts that are not attributable to changes in population size. For  
66 example, counts tend to be slightly lower than expected in an observer’s first year of service  
67 (Kendall et al. 1996) perhaps due to unfamiliarity with the route and survey. Within-observer  
68 effects have been used to investigate age-related change in counts (Link and Sauer 1998) and  
69 experimental changes in count protocols (Sauer et al. 2019).

70 Year effects reflect temporal and spatial changes in population change, the biologically relevant  
71 signals in noisy count data. Year effects are stratified by the intersection of states or provinces  
72 and Bird Conservation Regions (BCRs); the BCRs are physiographic regions that define  
73 major habitats relevant for birds across North America (Sauer et al 2013). While it is possible  
74 to estimate year effects at the stratum level without further modeling, estimation is greatly  
75 enhanced by the use of a hierarchical model in which year effects are treated as random

76 effects. The model currently used by the USGS (Model S) treats year effects as normally  
77 distributed with linearly trending means on the log scale; the slopes and intercepts of the  
78 regression coefficients, and the residual variance parameter are allowed to vary among strata,  
79 themselves as random effects.

80 Model S has the benefit of flexibility, being applicable to data of widely different quality. 493  
81 of the 548 species considered in this paper have data for at least 47 of the 50 years from  
82 1966-2015. Among these, sample sizes range from 113 to 99695, means from 0.05 to 83.5,  
83 coefficients of variation from 0.07 to 1.08, and numbers of strata ranging from 1 to 163. For  
84 many species, the data overwhelm the prior: the fitted trajectory (pattern of population  
85 of change) is distinctly nonlinear. For species with weak data, Model S is appropriate for  
86 estimation of a long term pattern of population trend.

87 Nevertheless, Model S is only one of many – infinitely many – models that can be fit to BBS  
88 data. The ease of fitting complex models via Markov chain Monte Carlo (MCMC) means that  
89 many models can be considered. Software for MCMC (e.g. JAGS, Plummer 2003; R2Jags, Su  
90 and Yajima 2015) requires nothing more of the analyst than the specification of models and  
91 priors; existing models are easily tweaked and new results produced. The merits of alternative  
92 models cannot be addressed on purely subjective grounds, such as interesting patterns in  
93 results, nor can a model be viewed as preferable because of its smaller standard errors,  
94 since the validity of the standard errors depends on the appropriateness of the model. Clearly,  
95 there is a need for objective model selection, based on sensibly defined criteria (Chatfield  
96 1995, Burnham and Anderson 2002, Link and Barker 2006, Hooten and Hobbs 2015).

97 This paper presents results of the first large-scale formal model selection exercise for the  
98 BBS. We compare 4 models for 548 species, using data for the period from 1966-2015.

<sub>99</sub> The four models are a  $2 \times 2$  cross-classification of models for year effects and models for  
<sub>100</sub> overdispersion. We compared models using the Bayesian predictive information criterion  
<sub>101</sub> (BPIC) and evaluated the Watanabe/Akaike information criterion (WAIC) as a convenient,  
<sub>102</sub> computationally fast alternative to BPIC (Gelman et al. 2014, Link and Sauer 2016b).

<sub>103</sub> We begin by describing a set of four candidate models for BBS data. These include the model  
<sub>104</sub> presently used in the USGS analyses, models with alternative structures for population change,  
<sub>105</sub> and models with alternative patterns of overdispersion relative to the Poisson distribution.  
<sub>106</sub> Next, we describe the BPIC and WAIC model selection criteria. Both relate to the posterior  
<sub>107</sub> predictive distribution, which is a well-known basis for evaluating model fit. We provide  
<sub>108</sub> details on the calculations involved in our model selection, then describe the methods we  
<sub>109</sub> use to compare results from the selected model with the model presently used by the USGS  
<sub>110</sub> (Model S). We report our primary results on model selection for 548 North American bird  
<sub>111</sub> species, and comment on the usefulness of WAIC as a surrogate for BPIC. Finally, we compare  
<sub>112</sub> results from the selected model with those from Model S.

<sub>113</sub> **Model set for BBS data**

<sub>114</sub> *Features shared among all models*

<sub>115</sub> The models we considered for BBS data are all overdispersed Poisson regressions of the form  
<sub>116</sub>  $Y_i|\lambda_i \sim P(\lambda_i)$ , where

$$\log(\lambda_i) = \Gamma_i + \Omega_i + \epsilon_i ; \quad (1)$$

<sub>117</sub> here  $\Gamma_i$ ,  $\Omega_i$ , and  $\epsilon_i$  are year, observer, and overdispersion effects, respectively.

<sub>118</sub> Year effects reflect the influence of bird abundance on counts. No attempt is made to model  
<sub>119</sub> absolute abundance because of the heroic assumptions needed to do so (Barker et al. 2018,  
<sub>120</sub> Link et al. 2018); instead,  $\Gamma_i$  is a measure of relative abundance, a function of parameters  
<sub>121</sub> describing spatial and temporal patterns in bird populations. In all of the models we consider,  
<sub>122</sub>  $\Gamma_i \equiv \gamma_{s(i),y(i)}$ ; here, indirect indexes  $s(i)$  and  $y(i)$  denote the stratum and year of the  $i^{th}$   
<sub>123</sub> observation. Thus year effects are described by a set of parameters  $\gamma_{s,y}$ , modeled as functions  
<sub>124</sub> of stratum  $s$  and year  $y$ .

<sub>125</sub> Observer effects reflect biologically irrelevant variation in counts related to differences among  
<sub>126</sub> observers, and differences within observers through time. In the models we consider, these  
<sub>127</sub> include a fixed effect  $\eta$  for an observer's first year of service on a route, and a mean-zero  
<sub>128</sub> normal random effect  $\omega_i$  for each combination of observer and route. Thus  $\Omega_i = \omega_{o(i)} + \eta f(i)$ ;  
<sub>129</sub> here,  $o(i)$  denotes the observer that produced count  $i$  and  $f(i)$  is an indicator variable for  
<sub>130</sub> whether count  $i$  is the observer's first count. Observer effects  $\omega_o$  are mean zero normal random  
<sub>131</sub> variables with precision (1/variance)  $\tau^\omega$ .

<sub>132</sub> Conditional on the mean parameter  $\lambda$ , the mean and variance of Poisson random variables  
<sub>133</sub> are identical. BBS counts are substantially more variable than their means would indicate,  
<sub>134</sub> hence overdispersion effects are modeled through the addition of a mean zero random effect  
<sub>135</sub>  $\epsilon_i$  in the linear predictor (1).

<sub>136</sub> *Slope, difference, and heavy-tailed models*

<sub>137</sub> The four models we consider are labeled D, DH, S, and SH. Model S is a slightly modified  
<sub>138</sub> version of the model that has been used for most BBS analyses since 2011 (e.g., Sauer et al.  
<sub>139</sub> 2014, Sauer and Link 2011). Labels S (for "slope"), D (for "difference") and H (for "heavy  
<sub>140</sub> tails") describe features of the models, which we now describe.

<sub>141</sub> The S models (S and SH) assume that year effects  $\gamma_{s,y}$  are conditionally independent and  
<sub>142</sub> normally distributed with precision  $\tau_s^\gamma$ ; the precision is allowed to vary among strata. The  
<sub>143</sub> expected value of  $\gamma_{s,y}$  is

$$E(\gamma_{s,y}) = S_s + \beta_s(y - y_0) ;$$

<sub>144</sub> here, the intercept  $S_s$  is a baseline abundance parameter for stratum  $s$ ,  $\beta_s$  is a trend parameter,  
<sub>145</sub> and  $y_0$  is a baseline year to center the regression. Baseline abundance and trend parameters  
<sub>146</sub> are allowed to vary by stratum;  $S_s$  and  $\beta_s$  are modeled as random effects, across strata.  
<sub>147</sub> The “slope” designation for models S and SH relates to the linear component of the model  
<sub>148</sub> described. This linear component is solely a prior expectation, in absence of data: actual year  
<sub>149</sub> effects vary, and the model is capable of detecting distinctly nonlinear population trajectories,  
<sub>150</sub> or patterns of population change.

<sub>151</sub> The D models (D and DH) replace the assumption that year effects  $\gamma_{s,y}$  are conditionally  
<sub>152</sub> independent, with the assumption that *differences* in year effects are conditionally independent.  
<sub>153</sub> Thus,  $\gamma_{s,y}$  is normally distributed with mean  $\gamma_{s,y-1}$  and precision  $\tau_s^\gamma$ . As with the slope  
<sub>154</sub> models, we fix a baseline year  $y_0$ , and set  $E(\gamma_{s,y_0}) = S_s$ , at the stratum mean.

<sub>155</sub> The D models have constant prior expectation, but like the S model, are capable of detecting  
<sub>156</sub> distinctly nonlinear population trajectories. With good data, estimated patterns of population  
<sub>157</sub> change for S and D models can be nearly identical (Link and Sauer 2016b). Differences  
<sub>158</sub> between fits relate to Bayesian shrinkage. Thinking of a graph of year effects against years, S  
<sub>159</sub> models shrink year effects vertically (toward a long term trend curve) while D models shrink  
<sub>160</sub> horizontally (towards values in adjacent years).

<sub>161</sub> In models S and D, extra-Poisson overdispersion effects  $\epsilon_i$  are assumed to be independent

and normally distributed with precision  $\tau_\epsilon$ . The normal model might not adequately account for extreme counts. Our experience is that extreme counts are a regular feature of BBS data, suggesting the need for a heavy-tailed alternative to the normal distribution. Thus the H models (SH and DH) specify a central t-distribution for  $\epsilon_i$  in modeling extra-Poisson variation, in place of a normal distribution. The t-distribution has scale parameter  $\tau^\epsilon$  and degrees of freedom parameter  $\nu$ . Following Juárez and Steel (2010) we use a Gamma distribution with mean 20 and variance 200 as an objective prior for  $\nu$ . In all four models, mean parameters are assigned flat normal priors, and scale parameters are assigned flat gamma priors. JAGS code for MCMC analysis is provided in Appendix S1.

### 171 *Composite trends*

Analyses of BBS data are typically summarized by estimates of trend and annual indices of abundance, computed for states, physiographic regions, countries, or for the entire survey area.

At the stratum level, annual indices have been based on expected counts in the region. These were of the form  $n_{s,y} = \exp(\gamma_{s,y} + \frac{1}{2}(\sigma_s^\gamma)^2 + \frac{1}{2}(\sigma_s^\epsilon)^2)$ , where  $\sigma_s^\gamma = 1/\sqrt{\tau^\gamma}$  and  $\sigma_s^\epsilon = 1/\sqrt{\tau^\epsilon}$  are the standard deviations of the random effects distributions of years effects and observer effects, respectively. These expected counts were derived under the assumption that that  $\epsilon_i$ 's and  $\omega_o$ 's follow normal distributions (Sauer and Link 2011); for the H models, these weights must be modified. We replace  $\sigma_s^\epsilon$  with a multiple of the t-distribution's scale parameter. Details are given in Appendix S2.

Regional trends and annual indices are derived statistics: the models are fit for all data for a species among the survey strata, statistics are computed at the level of survey strata and then aggregated among the survey strata to form regional estimates (e.g., Sauer and Link

185 2011). Composite annual indices for groups of strata are area-weighted stratum-level annual  
186 indices. Trend, defined as an interval-specific estimate of geometric mean yearly change, is  
187 computed as a ratio of annual indices for the last and first years of the interval, taken to the  
188 power  $1/(y_{\text{last}} - y_{\text{first}})$  (i.e., 1 over the length of the time interval, in years).

## 189 Model selection criteria

190 We begin by describing two model selection criteria, the BPIC and WAIC. Both will be seen  
191 to relate to the posterior predictive distribution, commonly used in model checking.

192 We denote the complete set of BBS counts for a given species, and associated covariates, by  
193  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_N)$  and  $\mathbf{X} = (X_1, X_2, \dots, X_N)$ , respectively. We assume a model  $M$  for the  
194 data, in which the observations  $Y_i$  are conditionally independent, given  $\mathbf{X}$ , with probabilities  
195  $\Pr(Y_i|X_i, M, \theta^M)$ .

196 Given a prior distribution for  $\theta^M$ , we can compute a posterior distribution  $\pi(\theta^M|\mathbf{D}, M)$ ,  
197 where  $\mathbf{D} = (\mathbf{Y}, \mathbf{X})$ . Using this, the posterior predictive distribution is calculated as

$$\text{ppd}(y|X, \mathbf{D}, M) = \int \Pr(y|X, M, \theta^M) \pi(\theta^M|\mathbf{D}, M) d\theta^M .$$

198 The distribution function  $\text{ppd}(y|X, \mathbf{D}, M)$  predicts the probability that a new observation  
199  $Y$  with covariate  $X$  will take the value  $y$ , given that  $Y$  is generated by model  $M$ , taking  
200 into account the uncertainty associated with the parameter  $\theta^M$ . The posterior predictive  
201 distribution is familiar as the basis of posterior predictive checks, used in assessing goodness  
202 of fit for Bayesian models.

203 Given a set of models  $\mathcal{M}$ , we seek the model  $M \in \mathcal{M}$  that is best able to predict new data  
204 based on the data at hand. A good predictive model  $M$  based on data  $\mathbf{D}$  should produce

205  $\text{ppd}(y|X, \mathbf{D}, M)$  that is close to the true but unknown (data generating) distribution, which  
 206 we denote by  $f(y|X)$ . Using the Kullback-Liebler divergence as a measure of closeness, it  
 207 turns out that we seek the model  $M \in \mathcal{M}$  that maximizes

$$\psi^M = \text{E}_f(\log(\text{ppd}(Y|X, \mathbf{D}, M))) \quad . \quad (2)$$

208 How are we to estimate this quantity? Given a hypothetical new observation  $\tilde{Y}$  sampled  
 209 from  $f(y|\tilde{X})$ ,

$$\log(\text{ppd}(\tilde{Y}|\tilde{X}, \mathbf{D}, M))$$

210 is an unbiased estimate of  $\psi^M$ . But what if we evaluate the log of the posterior predictive  
 211 distribution at the value of an observed datum  $(Y_i, X_i)$ ? The result,

$$\log(\text{ppd}(Y_i|X_i, \mathbf{D}, M)) \quad , \quad (3)$$

212 is an overestimate of  $\psi^M$ ; this happens because we have used observation  $Y_i$  in predicting  
 213 its own outcome. Leave-one-out cross-validation (LOOCV) mitigates this bias by instead  
 214 computing

$$B_i^M = \log(\text{ppd}(Y_i|X_i, \mathbf{D}_{-i}, M)) \quad ; \quad (4)$$

215 here,  $\mathbf{D}_{-i}$  is the data set  $\mathbf{D}$ , but with observation  $Y_i$  omitted. The necessity and benefit of  
 216 LOOCV are illustrated in Appendix S3. If using MCMC, calculation of (4) is accomplished  
 217 by averaging  $\Pr(Y_i|X_i, M, \theta^M)$  against draws of  $\theta_M$  from  $\pi(\theta^M|\mathbf{D}_{-i}, M)$ , then taking the

218 logarithm of the average.

219 The Bayesian predictive information criterion (BPIC) is defined as the sum across all  
220 observations of (4), viz.,

$$\text{BPIC} = \sum_{i=1}^N B_i^M = \sum_{i=1}^N \log(\text{ppd}(Y_i | X_i, \mathbf{D}_{-i}, M)) . \quad (5)$$

221 BPIC is regarded as the most reliable option currently available for selection among hierarchical  
222 models (Gelman et al. 2014, Link and Sauer 2016b).

223 Calculating BPIC is computationally intensive. For each of the  $N$  observations, calculating  
224 (4) involves calculation of  $\pi(\theta^M | \mathbf{D}_{-i}, M)$ , the posterior distribution for  $\theta^M$  based on the data  
225 set  $\mathbf{D}_{-i}$ . While importance sampling can be used to obtain samples of  $\pi(\theta^M | \mathbf{D}_{-i}, M)$  based  
226 on  $\pi(\theta^M | \mathbf{D}, M)$ , the procedure is unstable, and it is safest to perform a new MCMC analysis.

227 The computations required are substantial.

228 For example, the BBS White-winged Dove (*Zenaida asiatica*) data set has roughly the median  
229 number of observations, with 7389, collected on 24 strata. On a workstation equipped with  
230 Xeon E5-2630 v3 processors, MCMC with chain length 10000 takes approximately 12 minutes  
231 running on a single core. Running 7389 LOO analyses, even taking advantage of multicore  
232 processing on the system's 16 cores, would take over 3 days. And that's for a single model.

233 We estimate that the corresponding analysis for the widespread Mourning Dove (*Zenaida*  
234 *macroura*) with 98372 observations on 158 strata would take something over 3 months.

235 Computational expense can be reduced by computing BPIC on a subset of  $n \ll N$  observations  
236 (Link and Sauer, 2016b). We chose to conduct approximately  $n = 100$  analyses per species,  
237 rather than the average  $N = 18778$  analyses that would have been required for full computation

238 of BPIC.

239 The Watanabe/Akaike Information (WAIC) avoids the use of LOOCV, using (3) instead of  
240 (4) in estimating  $\psi^M$ . A bias correction term is added, so that WAIC is defined as

$$\text{WAIC} = \sum_{i=1}^N \{ \log(\text{ppd}(Y_i|X_i, \mathbf{D}, M)) - \text{Var}_{\text{post}}(\Pr(Y_i|X_i, M, \theta^M)) \} \quad , \quad (6)$$

241 where  $\text{Var}_{\text{post}}$  denotes the posterior variance of  $\theta^M$  based on the full data set  $\mathbf{D}$ . WAIC is  
242 asymptotically equivalent to BPIC (Gelman et al. 2014) and is much more easily computed  
243 than BPIC, because it requires only one analysis of the data. We calculated WAIC on the  
244 full sets of  $N$  observations, as well as on the subsets of  $n$  observations, to evaluate its use in  
245 approximating BPIC.

246 The BPIC and WAIC definitions given are totals across observations, with larger values  
247 indicating better predictive value of the data, based on the data  $\mathbf{D}$ . Some authors multiply  
248 by  $-2$  to put these criteria on the same (positive) scale as the Akaike, Bayes, and deviance  
249 information criteria, in which case smaller values are favored. The only scaling used in this  
250 paper will be by  $1/n$  to treat the subset BPIC criterion as an estimate of  $\psi^M$ . That is, we  
251 will write

$$\bar{B}^M = \frac{1}{n} \sum_{j=1}^n \log(\text{ppd}(Y_{i_j}|X_{i_j}, \mathbf{D}_{-i_j}, M)) \quad , \quad (7)$$

252 where  $\{i_1, i_2, \dots, i_n\}$  is a randomly selected subset of  $n$  indices from  $\{1, 2, \dots, N\}$ .

253 **Model selection calculations for BBS data**

254 Our goal was to select models for 548 species of birds surveyed by the BBS, with particular

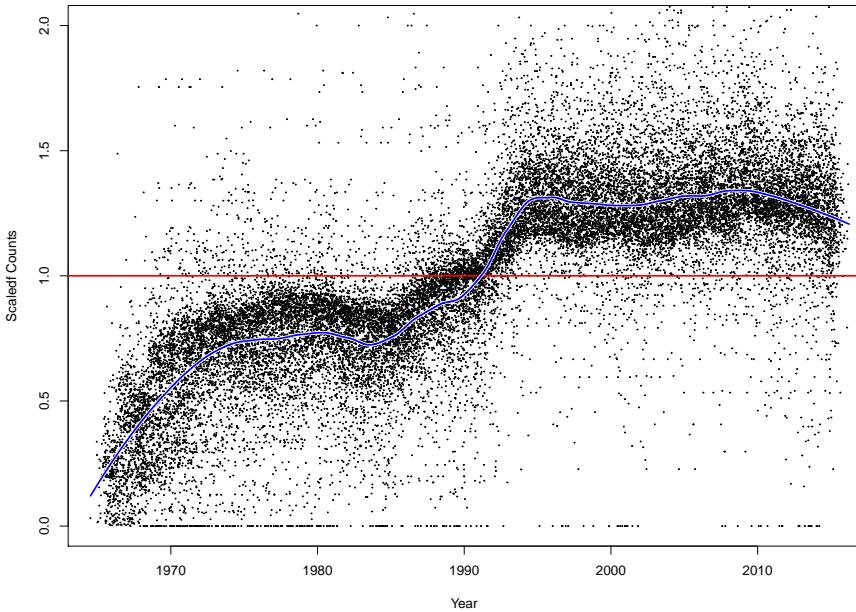


Figure 1: Scaled number of counts for the North American Breeding Bird Survey, each point corresponding to a species and a year. Each number has been scaled by the mean number of counts for the species over years subsequent to the species' first appearance in the data set. Blue curve is LOESS smooth. X-axis has been jittered to prevent overplotting.

255 interest in the trajectory component of the models. Our primary tool for model selection was  
 256 BPIC, calculated for a subset of roughly 100 counts per species. The geographic range of the  
 257 BBS and the density of routes have increased with time, meaning that species' data sets are  
 258 more heavily weighted toward recent years (Figure 1). We chose to sample 2 counts per year  
 259 for each species, so that years were equally represented in evaluating population trajectories.  
 260 Variation in the geographical extent of the survey resulted in different sample sizes, ranging  
 261 from 34 to 100, with 457 species having  $n \geq 96$ .  
 262 For each species and model, we began by running a Markov chain of length 10000, saving  
 263 the final state of the chain as a burned-in starting value for the subsequent  $n$  leave-one-out  
 264 analyses. The  $n$  LOO analyses are readily performed in parallel on a multicore processor.

265 Having selected  $n$  indices  $i_1, i_2, \dots, i_n$ , the  $j^{th}$  analysis uses the same MCMC code as the  
266 full analysis, but with observation  $i_j$  omitted, and a node calculating  $\Pr(Y_{i_j} | X_{i_j}, M, \theta^M)$   
267 monitored for a further 10000 Markov chain samples. The log of the posterior mean for this  
268 node is the  $i_j^{th}$  summand in the definition of BPIC, equation (5), which we denote by  $B_{i_j}$ .

269 Finally, for each species and model, we performed one more analysis of the full data set via  
270 MCMC (using the burned-in starting values previously) in order to compute WAIC values;  
271 memory requirements were substantial because of the need to monitor  $N$  nodes, one per  
272 observation. We found that chains of length 2500 were satisfactory for stable and precise  
273 estimates of the summands of WAIC (equation (6)); we denote the  $i^{th}$  value by  $W_i$ .

#### 274 Analysis of selected model

275 BBS trends are typically presented separately for “core” and “non-core” species (Sauer et al.  
276 2017). For 426 core species we present trends for the interval 1970-2015. Because BBS data  
277 covered limited geographic areas in the early years of the survey, we use 1970 as the first  
278 year for trend estimation. The 122 non-core species have limited BBS data prior to 1993; for  
279 these species we present trend results for the interval 1993-2015.

280 For the 548 species, we present trend results for Model S (our base model) and for the  
281 selected model based on the posterior probability that the model is best from the BPIC  
282 hierarchical model (i.e., Figure 2). For each model, we present the median estimated trend,  
283 2.5%CI, 97.5%CI, and the annual index for the midyear. Trend estimates represent yearly  
284 percentage change for the interval 1970-2015 (for core species) or 1993-2015 (for non-core  
285 species). We also present the sample size, the number of survey routes used in the analysis  
286 for the species. For these trend estimates, we summarize differences in results between model  
287 S to the selected model.

288 If the model selection procedure tends to generally favor D models or S models, the model  
289 selection procedure might result in systematic differences in trends between results from the  
290 default model and selected model. Prior work has shown that the differences in year effects  
291 parameterizations of D vs S models can lead to systematic differences in trends. S models  
292 model a slope parameter, on which year effects are deviations, while the D models directly  
293 model year to year changes. Because the BBS had very limited data in the early years of the  
294 survey, we expect that these differences between the models lead to more extreme trends for  
295 the S models (as the trajectory in the early years of the survey are dominated by the slope  
296 parameter in the model). We also predict that D models will be less precise for years with  
297 weak data, for similar reasons, i.e., the slope-based model is dominated by a slope parameter  
298 while the D models are dominated by poorly estimated year effects.

299 We used paired t-tests to determine whether systematic differences in trends existed between  
300 trend results from the default model and the selected model at the continental (survey-wide)  
301 scale of summary. We also evaluated whether half-widths of 95% credible intervals of trend  
302 were consistently different between default model and the selected models. We also identify  
303 species for which the 95% CIs of trend estimates did not overlap between the default model  
304 and the selected model. Although this criterion is of limited quantitative significance due to  
305 lack of independence, it does highlight species for which the change in model selected can have  
306 large effects on the trend estimates. For these species, we provide graphs of annual indices  
307 for default models and the selected models. We conducted summary analyses separately for  
308 core and non-core species.

309 **Results: Model selection using BPIC**

310 *Basic BPIC results*

311 Let  $\bar{B}_s^M$  denote the BPIC value based on the subsample of size  $n \ll N$  for species  $s =$   
 312  $1, 2, \dots, 548$ , as defined at equation (7). Ranking models based simply on raw values  $\bar{B}_s^M$ ,  
 313 model D is ranked as best 141/548 times (25.7%), model DH 203 times (37%), model S 81  
 314 times (14.8%), and model SH 123 times (22.4%). On the other hand, model D is ranked as  
 315 worst 138 times (25.2%), model DH 70 times (12.8%), model S 213/548 times (38.9%), and  
 316 model SH 127 times (23.2%). This quick appraisal (summarized in Table 1) suggests that  
 317 across species, there is a tendency for DH to be favored and S to be disfavored.

	D	DH	S	SH
<b>Ranked Best</b>	141	203	81	123
	25.7%	37%	14.8%	22.4%
<b>Ranked Worst</b>	138	70	213	127
	25.2%	12.8%	38.9%	23.2%

Table 1: Rank frequencies (based on sampled BPIC values) across species

318 Associated with each vector  $\hat{\mathbf{F}}_s = (\bar{B}_s^D, \bar{B}_s^{DH}, \bar{B}_s^S, \bar{B}_s^{SH})'$  is an estimated covariance matrix  
 319  $\hat{\Sigma}_s$  (the sample covariance matrix for the  $n$  sampled values) from which one can calculate  
 320 standard errors of differences, such as  $(\bar{B}_s^D - \bar{B}_s^{DH})$ . Using these, we calculate z-statistics to  
 321 test null hypotheses of equal support for models (see Discussion in Link and Sauer 2016b, pg  
 322 1756). Most of the differences are not precisely estimated, with the result that relatively few  
 323 of the model comparisons are significant at  $\alpha = 0.05$  (274 out of  $548 \times 6$ , or 8.3%).  
 324 Nevertheless, inspection of Table 2 suggests that general tendencies emerge in comparisons of  
 325 models, aggregated across species. The last row of Table 2 (labeled ‘Null’) gives the expected  
 326 frequencies for a sample of 548 normal random variables; these are the expected frequencies  
 327 under a null hypothesis of no difference in model fits. For each of the 6 model comparisons,  
 328 the distribution of results is shifted to the right or left relative to this baseline. Comparisons  
 329 ‘D vs DH’ and ‘S vs SH’ have outcome frequencies shifted to the left, favoring DH over D, and

	$< -1.96$	$[-1.96, 0)$	$[0, 1.96)$	$\geq 1.96$	Mean	SD
D vs DH	21	315	203	9	-0.28	1.03
D vs S	6	205	289	48	0.37	1.13
D vs SH	17	240	263	28	0.07	1.12
DH vs S	5	176	314	53	0.51	1.06
DH vs SH	4	196	309	39	0.39	1.08
S vs SH	36	299	205	8	-0.32	1.09
Null	14	260	260	14	0	1.00

Table 2: Paired z-test statistics for BPIC differences.

330 SH over S. The other 4 comparisons have outcome frequencies shifted to the right, favoring

331 D over S and SH, and DH over SH.

332 These results are consistent with an ordering  $DH > D > SH > S$  of general tendencies in

333 preferences among models (here “ $>$ ” means “has better predictive value than”). Heavy tailed

334 models and difference models are preferred, and of the two innovations, difference models

335 offer the greater gain.

336 We will say that model  $M_1$  is favored over model  $M_2$  if the z-statistic is greater than 1.96.

337 There are only 3 species for which DH is favored over all, and only 1 where D is favored over

338 all. S and SH are never favored over all.

339 For 33 species, model D or DH is favored over both S and SH. Of these, there is no clear

340 preference for D or DH in 29 cases, with DH being favored over D in 3 cases, and D over DH

341 in 1. On the other hand, there are only 3 species for which S or SH is favored over both D

342 and DH; in none of these is either S or SH favored.

343 For 15 species, heavy-tailed models DH or SH are favored over the non-heavy-tailed alternatives

344 D and S. For 3 of these, DH is favored over SH; in none is SH favored over DH. On the

345 other hand, there are only 3 species for which D or S is favored over both DH and SH. D is

346 preferred over S in one of these; neither S nor D is favored in the others.

347 *Species treated as replicates in hierarchical model for BPIC*

348 Given enough data, the BPIC criterion will distinguish between the fits of alternative models.

349 Thus the relatively small number of species for which clear preferences exist is due to the

350 limited subsampling of counts. The tendency for first difference models D and DH to be

351 favored over models S and SH, and for heavy-tailed models DH and SH to be favored over

352 D and S is appropriately conducted using a hierarchical model, across species. Hierarchical

353 modeling also enhances model selection for individual species by considering them collectively,

354 “borrowing strength from the ensemble” (Morris 1983, Louis 1984, Rubin 1984).

355 For species  $s$ , let  $\hat{\Delta}_s = C\hat{F}_s$ , where  $C$  is the  $3 \times 4$  contrast matrix

$$C = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix} .$$

356 That is,  $\hat{\Delta}_s = (\bar{B}_s^D - \bar{B}_s^{DH}, \bar{B}_s^S - \bar{B}_s^{DH}, \bar{B}_s^{SH} - \bar{B}_s^{DH})'$ . We treat  $\hat{\Delta}_s$  as having a trivariate

357 normal distribution with mean vector  $\Delta_s$  and known covariance matrix  $V_s = C\hat{\Sigma}_s C'$ . We

358 assume a trivariate normal distribution for the means, i.e.,  $\Delta_s \sim \mathcal{N}_3(\mu, W)$ , assign a flat

359 normal prior to  $\mu$ , and a vague inverse Wishart prior to  $W$ . We analyzed this hierarchical

360 model for BPIC values, obtaining posterior distributions for  $\Delta_s$  given the complete collection

361 of estimates,  $\hat{\Delta}_\cdot = (\hat{\Delta}_1, \hat{\Delta}_2, \dots, \hat{\Delta}_{548})'$  as data.

362 The latent vector  $\Delta_s$  consists of differences  $\psi_s^D - \psi_s^{DH}$ ,  $\psi_s^S - \psi_s^{DH}$  and  $\psi_s^{SH} - \psi_s^{DH}$ , with

363  $\psi^M$  defined at equation (2) and subscript  $s$  added for species. Recalling that larger values

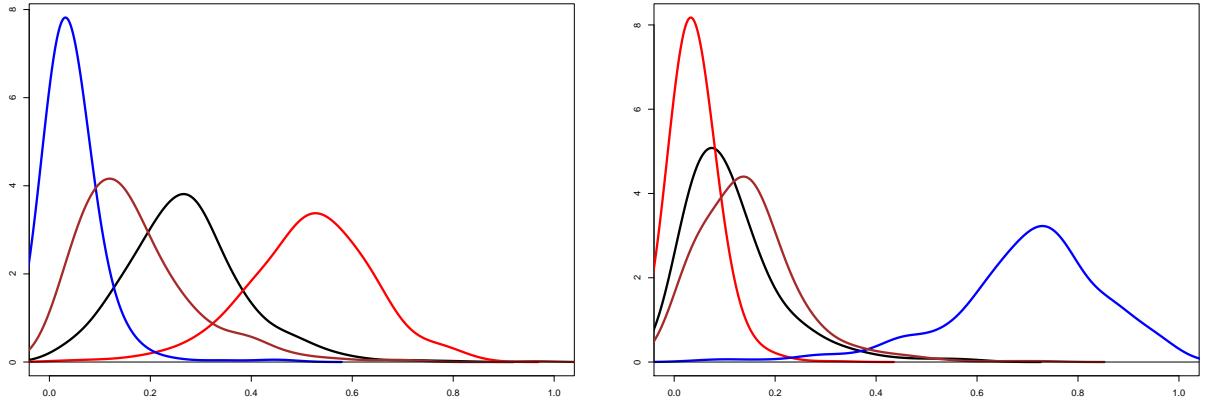


Figure 2: Smoothed density of posterior probability model is best (left panel,  $\pi_s^M$ ) and worst (right panel,  $\rho_s^M$ ) across 548 BBS species, for models M = D (black), DH (red), S (blue) and SH (brown).

of  $\psi^M$  indicate better expected predictive value for the combination of data and model,  
 $\Delta_{s,1} > 0$  means that model D is favored over model DH,  $\Delta_{s,2} > 0$  means that model S is  
favored over model DH, and  $\Delta_{s,3} > 0$  means that model SH is favored over model D. The  
four models can be ranked based on  $\Delta_s$ ; for example,  $\Delta_s = (-1.2, 3.8, -2.5)$  implies the  
ordering S > DH > D > SH. Thus we obtained values  $\pi_s^M = \Pr(M \text{ is best for } s | \hat{\Delta}_s)$  and  
 $\rho_s^M = \Pr(M \text{ is worst for } s | \hat{\Delta}_s)$ . Results are summarized in Table 3 and Figure 2; results by  
species are given in Appendix S4.

*Results: WAIC as a surrogate for BPIC*

	D	DH	S	SH
Highest $\pi_s^M$	53	458	4	33
Mean $\pi_s^M$	27%	51.5%	4.6%	17%
Highest $\rho_s^M$	13	0	524	11
Mean $\rho_s^M$	11.5%	4.1%	69.9%	14.6%

Table 3: Posterior probabilities that model is best ( $\pi_s^M$ ) or worst ( $\rho_s^M$ ), summarized across species.

372 In an earlier paper investigating a set of 4 models in application to 20 BBS species (Link  
373 et al. 2017) we noted that the components of WAIC ( $W_i$ ) were poor surrogates for the  
374 components of BPIC ( $B_i$ ). We confirmed that finding in the present study. Figure 3 plots  
375 values  $W_i$  against corresponding values  $B_i$  for model D, using data for all 548 species. The  
376 pattern is evident at the species level, and across models. Gray lines indicate the 1<sup>st</sup>, 5<sup>th</sup>,  
377 25<sup>th</sup>, and 50<sup>th</sup> percentiles of BPIC values. Note that large values of  $B_i$  (suggesting good  
378 fit of the model) tend to be matched well by large values of  $W_i$ . However, small values of  
379  $B_i$  (e.g. below the 5<sup>th</sup> percentile) have corresponding values of  $W_i$  that are too large. The  
380 corresponding observations are those corresponding to poorest model fit, and which make the  
381 largest contributions to evaluating the model. Thus WAIC appears to overstate the predictive  
382 value of the model and data.

383 We considered using the  $n$  calculated values of  $B_i$  with the corresponding values of  $W_i$  in order  
384 to establish a regression-based prediction of  $B_i$  for the  $N - n$  values not directly calculated,  
385 based on the easily obtained corresponding values of  $W_i$ . Unfortunately, the resulting  
386 prediction variances were large enough to render the predictions unreliable. Additionally, we  
387 could not establish a consistent regression relation between  $W_i$ 's and  $B_i$ 's across species.

388 Nevertheless, the relative ease with which WAIC can be calculated, and the potential  
389 advantages of using all  $N$  observations rather than the subset of  $n$  used in our BPIC analyses,  
390 invite further comparison. For comparability with our BPIC sampling of equal numbers of  
391 observations per year, we computed a weighted version of WAIC. This weighted WAIC is of  
392 the same form as (6), but includes weights  $w_i \propto 1/n_{y(i)}$ , where  $n_y$  is the number of counts in  
393 year  $y$ , and  $y(i)$  is the year in which the  $i^{th}$  count was conducted.

394 We compare model rankings for the weighted WAIC to the rankings by posterior probability

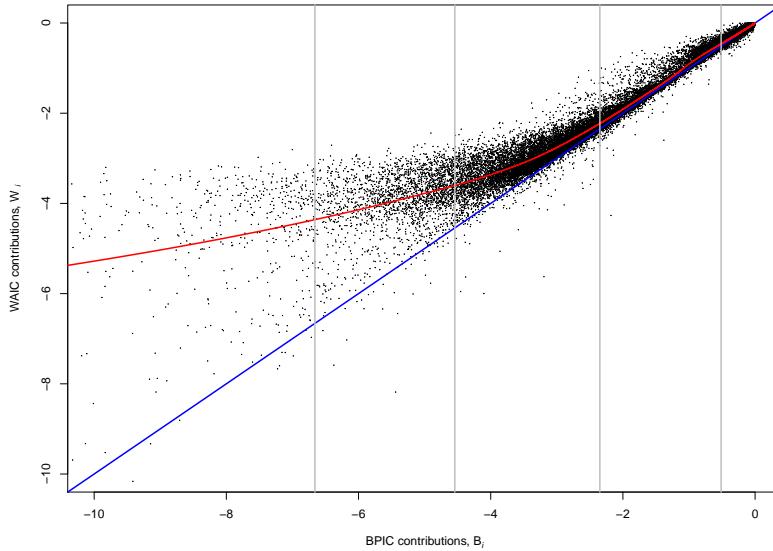


Figure 3: Components of WAIC for Model D ( $W_i^D$ ) plotted against corresponding components of BPIC ( $B_i^D$ ), 548 species. Blue line is identity; red is LOESS curve. Gray lines indicate 1<sup>st</sup>, 5<sup>th</sup>, 25<sup>th</sup>, and 50<sup>th</sup> percentiles of  $B_i^D$  values.

395 in the hierarchical analysis of sampled BPIC values. Across the 548 species, the weighted  
 396 WAIC assigns rank 1,2,3,4 to the BPIC best model at rates of 46.2%, 24.8%, 15.0%, and  
 397 14.0%, respectively.

398 Figure 4 plots the results of a logistic regression of the success rate of the weighted WAIC in  
 399 identifying the BPIC model, conducted to examine how the success rate increases with evidence  
 400 in favor of the best model. That is, we let  $S_i$  be the indicator for species  $i = 1, 2, \dots, 548$ ,  
 401 and modeled

$$\text{logit}(\Pr(S_i = 1)) = \alpha + \beta \text{ logit}(\pi_i) ,$$

402 where  $\pi_i$  is the posterior probability of the best supported model for species  $i$  (note that with  
 403 4 models in the model set, this must be at least 0.25).

404 Finally, we considered the relation between complete set of rankings from the weighted WAIC

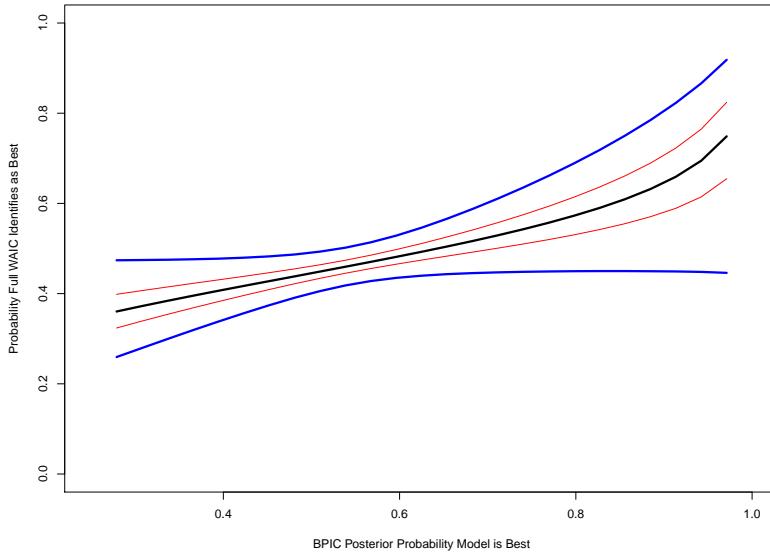


Figure 4: WAIC success rate in identifying BPIC best model, as function of posterior probability of BPIC best model. Red curves are 25<sup>th</sup> and 75<sup>th</sup> percentiles of estimation uncertainty in logistic regression; blue curves are corresponding 2.5<sup>th</sup> and 97.5<sup>th</sup> percentiles.

405 and the BPIC. Given BPIC ranks {1,2,3,4}, the corresponding WAIC rankings were {1,2,3,4},  
 406 {2,1,3,4}, {1,3,2,4}, {1,2,4,3}, or {2,1,4,3} for 50.2% of the species. Simply guessing, the  
 407 chance of doing that well is 20.8%.

408 Taken together, these results show that the WAIC criterion is a weak surrogate for BPIC.

409 **Results: Selected model versus traditional model**

410 Results for the 548 species are presented in Appendix S4. We note that 4 species had Model  
 411 S as the preferred model, and we include those results in summary statistics. These species  
 412 cover a large array of sampling and life history situations, varying in samples from 3 to  
 413 4423 routes (median=341), and in estimated abundance from 0.003 to 20073 (median=48.9).  
 414 Because D models were selected for the majority of species, we expected to see systematic  
 415 differences in results between default model and the selected models.

<sup>416</sup> For 426 core species with trends estimated for 1970-2015, we see no consistent differences  
<sup>417</sup> between trends from the default model and the selected model (mean difference, default –  
<sup>418</sup> selected = 0.04%/yr, 95%CI: -0.073, 0.016); for the 122 non-core species the default models  
<sup>419</sup> were consistently higher (mean difference = 0.99 %/yr, 95%CI: 0.024, 1.740), presumably  
<sup>420</sup> reflecting the effects of the differences between S and D models in regions with limited data.

<sup>421</sup> Comparisons of half-widths of credible intervals indicated that trends tended to be less precise  
<sup>422</sup> in the selected model than in the default model for the core species results (mean difference  
<sup>423</sup> = -0.40, 95%CI: -0.558, -0.236), but for the non-core species results the half widths of the  
<sup>424</sup> CIs were greater for the default model than for the selected models (mean difference = 2.96,  
<sup>425</sup> 95%CI: 2.180, 3.733).

<sup>426</sup> Eighteen of the core species had trend estimates for which the 95% CIs did not overlap  
<sup>427</sup> between default model and the selected model. We present the annual indices for these  
<sup>428</sup> species in Appendix S5.

## <sup>429</sup> Discussion

<sup>430</sup> Choosing among alternative statistical models for analysis of BBS data is an important means  
<sup>431</sup> of advancing our understanding of population change in North American birds. Choice of  
<sup>432</sup> appropriate models is complicated by the size of the dataset, the unruly distributional nature  
<sup>433</sup> of the data, the complexity of the models, and the temporal and spatial extent of the survey.

<sup>434</sup> Leave-one-out-cross-validation provides a reasonable approach for model selection, and the  
<sup>435</sup> BPIC statistic, as an observation-based metric of fit, provides greats flexibility for evaluating  
<sup>436</sup> temporal and spatial patterns of model fit. The challenge has been to apply the approach in  
<sup>437</sup> a computationally-feasible manner to the 548 BBS species for which analyses are conducted.

<sup>438</sup> Here, we have sampled observations for BPIC in a temporally-balanced design to assess

439 model fit for a model set that contrasts 2 alternative parameterizations of year effects and 2  
440 overdispersion distributions, and apply a hierarchical model to the model selection results to  
441 compute posterior probabilities models are best for species. Results suggest that models that  
442 use the D year effects parameterization and t-distribution based overdispersion are favored  
443 for most species.

444 Our results begin with simple comparisons of sampled BPIC values, then follow with statistical  
445 tests comparing fits between models, ending with a hierarchical model that allows us to  
446 obtain posterior probabilities that each model is best for a given species. Each successive  
447 refinement of analysis enhanced the view of DH and D as the top-ranked models. Based  
448 on these results, it seems sensible to abandon the S model as our default model as it has  
449 the highest posterior probability that it is the worst model for 524 species and the highest  
450 posterior probability that it is the best model for only 4 species. Model DH seems a much  
451 better candidate as a default model, with a highest posterior probability that it is the worst  
452 and best model for 0 and 458 species, respectively.

453 A concern with use of BPIC is the challenge to fully implement it for most BBS species. Our  
454 subsample analysis, based on  $n \ll N$  observations, was necessitated by the computational  
455 burden of computing values  $B_i^M$  (equation 4), the components of BPIC. As it turns out,  
456 the variation in  $B_i^M$  across models is small relative to the variation across observation  $Y_i$ .  
457 Sampling a larger portion of the observations would clearly lead to increased power of  
458 statistical tests among models and also provide more species-specific information to inform  
459 the hierarchical modeling. For species where estimates of population change differ among  
460 models, it would seem prudent to conduct larger sampling of observations to better estimate  
461 total BPIC.

<sup>462</sup> The numbers of BBS counts have changed over time, as has the spatial extent of the survey.  
<sup>463</sup> These features of the BBS need to be accommodated in model selection. In the interest of  
<sup>464</sup> optimizing selection of temporal pattern, we chose to balance years in our sampling of BPICs.  
<sup>465</sup> We suspect that if we had not included the temporal balancing of sampling, the S models  
<sup>466</sup> would have been more favored, as presumably a completely random sample of observations  
<sup>467</sup> would have placed greater emphasis on years in the middle and later years where more routes  
<sup>468</sup> were sampled.

<sup>469</sup> Larger samples of BPIC components  $B_i^M$  would be useful in evaluating model fit in specific  
<sup>470</sup> areas and periods of interest. For example, species such as Barn Swallow (*Hirundo rustica*)  
<sup>471</sup> show large differences in population trajectories in the early years of the survey (Appendix  
<sup>472</sup> S5) between S and D models, and additional evidence of the superiority of a D model for  
<sup>473</sup> these years could be provided by sampling additional observations from these years. Barn  
<sup>474</sup> Swallows also show striking regional differences in population trajectories, with increases in  
<sup>475</sup> the southern United States countered by strong declines in the central and northern United  
<sup>476</sup> States and Canada (Sauer et al. 2017). For this species, it would be of great interest to  
<sup>477</sup> evaluate both spatial and temporal patterns of model fit by mapping observation-specific  
<sup>478</sup> values of  $B_i^M$ .

<sup>479</sup> WAIC offers a potential solution to these concerns about small sample sizes, as it is easily  
<sup>480</sup> computed for all observations. Unfortunately, results presented here confirm concerns about  
<sup>481</sup> WAIC as a poor surrogate for BPIC. Matched against the yardstick of BPIC, WAIC results  
<sup>482</sup> appear to be somewhat better than chance, but only match the BPIC rankings 21% of the  
<sup>483</sup> time. In particular, the differential bias in WAIC with magnitude of BPIC further limits the  
<sup>484</sup> value of WAIC in evaluating temporal and spatial patterns of model fit.

485 The BPIC values presented here have informed our determination of the relative merits of  
486 the models in our model set, but were computed at significant cost in terms of computational  
487 effort, data storage, and personnel time. Hopefully, they also have some value in helping to  
488 design future model selection activities with different model sets. In the short term, for studies  
489 that propose alternative models (such as semiparametric smooths to model year effects, or  
490 the value of covariates to examine phenology effects on counts) it would be possible to use  
491 the data structure of years and sampled observations from the current study and estimate  
492 BPIC components  $B_i^M$  for the new models, based on the same samples of observations. BPIC  
493 for the new models could then be compared directly with the existing data.

494 Model selection within the model set described here allows for 2 generalizations with regard  
495 to model structure: (1) the D models, in which year effects are defined in terms of changes  
496 from adjacent years, seems generally preferable to the S models, in which year effects are  
497 defined as deviations from a consistent underlying slope parameter, and (2) the models that  
498 allow for t-distributed overdispersion are generally to be preferred over models that allow for  
499 normally distributed overdispersion. However, empirical results indicate that choice of model  
500 does not result in dramatically different views of population change within species. This is  
501 comforting, as a strong dependence on models (that are admittedly sometimes difficult to  
502 discriminate using our model selection tools) tends to undermine the credibility of results.  
503 Also, the distinction between D and S models can be thought of as a difference in process  
504 priors for model trajectories. For almost all species the large sample sizes appear adequate  
505 for the data to overwhelm the prior.

506 However, as illustrated by our evaluation of species with non-overlapping credible intervals of  
507 selected and default trend estimates, differences tend to occur when a pair of conditions exists:

508 (1) the species is undergoing dramatic and consistent population changes, and (2) the species  
509 has very limited data in the early years of the survey. This combination of circumstances  
510 leads to cases in which the slope parameter in the S models is the dominant expression of  
511 population change and the year effects are poorly estimated. Particularly in periods with  
512 limited data, predictions of annual indices based on the slope parameter are not particularly  
513 informative as they only lead to a linear (on a log scale) prediction of change. On the other  
514 hand, the D model does not have the constraint of a consistent prediction of change but  
515 instead only predicts change based on the information in the interval and a prior expectation  
516 of no change. Thus, the absence of information in the S model predicts the yearly change  
517 based on the prior defined by the slope parameter while the D model has a prior expectation  
518 of zero change.

519 The preeminence of the H models also conforms to our expectations regarding the BBS  
520 dataset. In our experience, BBS data tend to have many extreme observations. These extreme  
521 observations lead to lack of fit; the t-distributions modeling overdispersion in the H models  
522 accommodate these extreme observations. However, the H models introduce a complication  
523 in that the mathematical expectation of an exponentiated t-value is infinite; we must use  
524 an alternative characterization of a typical count. In Appendix S2, we provide both a clear  
525 quantitative rationale for estimation of expected counts in the case of normal overdispersion  
526 and suggest an analogue for use in the case of t-distributed overdispersion. To our knowledge  
527 this is the first published rationale for these abundance indices.

## 528 Acknowledgments

529 We thank the thousands of volunteers who have contributed to the North American Breeding  
530 Bird Survey. Any use of trade, product, or firm names in this publication is for descriptive

531 purposes only and does not imply endorsement by the U.S. Government.

532 **Literature cited**

533 Barker, R.J., M.R. Schofield, W.A. Link, and J.R. Sauer. 2018. On the reliability of N-mixture  
534 models for count data. *Biometrics* 74(1):369–377.

535 Burnham, K.P., and D.R. Anderson. 2002. Model selection and multimodel inference: a  
536 practical information-theoretic approach. Second edition. Springer-Verlag, New York,  
537 New York, USA.

538 Chatfield, C. 1995. Model uncertainty, data mining and statistical inference (with discussion).  
539 *Journal of the Royal Statistical Society (London), Series A* 158:419–466.

540 Gelman, A., J. Hwang, and A. Vehtari. 2014. Understanding predictive information criteria  
541 for Bayesian models. *Statistics and Computing* 24:997–1016.

542 Hooten, M.B., and N.T. Hobbs. 2015. A guide to Bayesian model selection for ecologists.  
543 *Ecological Monographs* 85:3–28.

544 Juárez, M.A. and M.F. Steel. 2010. Model-based clustering of non-Gaussian panel data based  
545 on skew-t distributions. *Journal of Business and Economic Statistics* 28:52–66.

546 Kendall, W.L., B.G. Peterjohn, and J.R. Sauer. 1996. First-time observer effects in the  
547 North American Breeding Bird Survey. *The Auk* 113:823–829.

548 Link, W.A., and R.J. Barker. 2006. Model weights and the foundations of multimodel  
549 inference. *Ecology* 87:2626–2635.

550 Link, W.A., J.R. Sauer, and D.K. Niven. 2017. Model selection for the North American  
551 Breeding Bird Survey: A comparison of methods. *Condor* 119(3):546–556.

- 552 Link, W.A., M.R. Schofield, R.J. Barker, and J.R. Sauer. 2018. On the robustness of  
553 N-mixture models. *Ecology* 99:1547–1551.
- 554 Link, W.A., and J.R. Sauer. 1998. Estimating relative abundance from count data. *Austrian  
555 Journal of Statistics*, 27:83–97.
- 556 Link, W.A., and J.R. Sauer. 2016a. Modeling participation duration, with application to  
557 the North American Breeding Bird Survey. *Communications in Statistics – Theory and  
558 Methods* 45:6311–6320.
- 559 Link, W.A., and J.R. Sauer. 2016b. Bayesian Cross-Validation for Model Evaluation and  
560 Selection, with Application to the North American Breeding Survey. *Ecology* 97:1746–  
561 1758.
- 562 Louis, T.A. 1984. Estimating a population of parameter values using Bayes and empirical  
563 Bayes methods. *Journal of the American Statistical Association* 79:393–398.
- 564 Morris, C.N. 1983. Parametric empirical Bayes inference: theory and applications. *Journal  
565 of the American Statistical Association* 78:47–55.
- 566 North American Bird Conservation Initiative. 2016. The State of North America’s  
567 Birds 2016. Environment and Climate Change Canada: Ottawa, Ontario. 8 pages.  
568 [www.stateofthebirds.org](http://www.stateofthebirds.org).
- 569 Plummer, M. 2003. JAGS: A program for analysis of Bayesian graphical models using Gibbs  
570 sampling. In *Proceedings of the 3rd International Workshop on Distributed Statistical  
571 Computing* (DSC 2003) March 20-22, Vienna, Austria (K. Hornik, F. Leisch, and A.  
572 Zeileis, Editors). pp. 1–10.
- 573 Rubin, D.B. 1984. Bayesianly justifiable and relevant frequency calculations for the applied

- 574 statistician. *Annals of Statistics* 12:1151–1172.
- 575 Rosenberg, K.V., J.A. Kennedy, R. Dettmers, R.P. Ford, D. Reynolds, J.D. Alexander, C.J.  
576 Beardmore, P.J. Blancher, R.E. Bogart, G.S. Butcher, A.F. Camfield, A. Couturier, D.  
577 W. Demarest, W.E. Easton, J.J. Giocomo, R.H. Keller, A.E. Mini, A.O. Panjabi, D.N.  
578 Pashley, T.D. Rich, J.M. Ruth, H. Stabins, J. Stanton, T. Will. 2016. Partners in Flight  
579 Landbird Conservation Plan: 2016 Revision for Canada and Continental United States.  
580 Partners in Flight Science Committee. 119 pp.
- 581 Sauer, J.R., and W.A. Link. 2011. Analysis of The North American Breeding Bird Survey  
582 Using Hierarchical Models. *The Auk* 128:87–98. +Supplemental Online Material.
- 583 Sauer J.R., W.A. Link, J.E. Fallon, K.L. Pardieck, and D.J. Ziolkowski Jr. 2013. The North  
584 American Breeding Bird Survey 1966–2011: Summary Analysis and Species Accounts.  
585 North American Fauna 79:1–32. + 900pp Supplemental material.
- 586 Sauer, J.R., J.E. Hines, J.E. Fallon, K.L. Pardieck, D.J. Ziolkowski Jr., and W.A. Link. 2014.  
587 The North American Breeding Bird Survey, Results and Analysis 1966 – 2013. Version  
588 01.30.2015 USGS Patuxent Wildlife Research Center, Laurel, MD.
- 589 Sauer, J.R., D.K. Niven, K.L. Pardieck, D.J. Ziolkowski Jr, and W.A. Link. 2017. Expanding  
590 the North American Breeding Bird Survey analysis to include additional species and  
591 regions. *Journal of Fish and Wildlife Management* 8(1):154–172.
- 592 Sauer, J.R., W.A. Link, D.J. Ziolkowski, K.L. Pardieck, and D.J. Twedt. 2019. Consistency  
593 counts: Modeling the effects of a change in protocol on Breeding Bird Survey counts. *The  
594 Condor* 121: duz009.
- 595 Su, Y.-S., and Yajima M. 2015. R2jags: Using R to run ‘JAGS.’ Version 0.5-7.

<https://CRAN.R-project.org/package/R2jags>

597 **Supporting Information:** William A. Link, John R. Sauer, Daniel K. Niven. Model  
598 Selection for the North American Breeding Bird Survey with observations on the WAIC and  
599 BPIC model selection criteria. Ecology.

600 **Appendix S1: JAGS code for models**

601 The four models we consider are labeled D, DH, S, and SH. Labels S (for “slope”), D (for  
602 “difference”) and H (for “heavy tails”) describe features of the models. A JAGS model  
603 statement for S is given in Figure 5. We use this to describe common features of the models,  
604 then in Figure 6 present changes to the year effects needed for D models. Changes needed for  
605 the H models are described in the text of this Appendix.

606 The `Data model` portion of the JAGS code is the same for all 4 models. Count  $i$  is modeled  
607 as a Poisson random variable with mean  $\lambda_i$ . We suppose that  $\log(\lambda_i)$  is a random variable  
608 with mean

$$E(\log(\lambda_i)) = \gamma_{s(i),y(i)} + \omega_{o(i)} + f(i) \eta \quad ;$$

609 here  $\gamma_{s,y}$  is the expected count in stratum  $s$  at year  $y$ ,  $\omega_o$  is the effect of observer  $o$ , and  $\eta$  is  
610 an offset for the effect of an observer’s first year of service (indicated by a zero-one variable  
611  $f(i)$ ). Note the use of indirect indexing:  $s(i)$ ,  $y(i)$ , and  $o(i)$  index the stratum, year, and  
612 observer corresponding to the  $i^{th}$  observation. Observer effects are treated as a sample of a  
613 mean zero normal distribution with precision  $\tau_\omega$ .

614 In models S and D,  $\log(\lambda_i)$  is assumed to be normally distributed with precision (1/variance)  
615  $\tau_\epsilon$ , reflecting extra-Poisson variation. Heavy-tailed alternative models SH and DH assume  
616 a t-distribution in place of a normal distribution in modeling extra-Poisson variation. This  
617 entails two changes to the code in Figure 5. First, we change the definition of `loglambda[i]`

```

# JAGS model statement for BBS data #
model=function(){

  # Data model #
  for(i in 1:ncounts){
    E.loglambda[i] <- gamma[year[i],strat[i]]+omega[obser[i]]+eta*firstyr[i]
    loglambda[i] ~ dnorm(E.loglambda[i],tau.epsilon)
    log(lambda[i]) <- loglambda[i]
    count[i] ~ dpois(lambda[i])
  }
  tau.epsilon ~ dgamma(0.001,0.001)

  # Observer effects #
  for(o in 1:nobservers){
    omega[o] ~ dnorm(0.0,tau.omega)
  }
  tau.omega ~ dgamma(0.001,0.001)
  eta ~ dnorm(0.0,1.0E-6)

  # Year effects #
  for(s in 1:nstrata){
    for(y in 1:nyears){
      E.gamma[s,y] <- stratum[s]+beta[s]*(y-y.baseline)
      gamma[s,y] ~ dnorm(E.gamma[s,y],tau.gamma[s])
    }
    beta[s] ~ dnorm(mu.beta,tau.beta)
    stratum[s] ~ dnorm(mu.strata,tau.strata)
    tau.gamma[s] ~ dgamma(0.001,0.001)
  }
  mu.beta ~ dnorm(0.0,1.0E-6)
  tau.beta ~ dgamma(0.001,0.001)
  mu.strata ~ dnorm(0.0,1.0E-6)
  tau.strata ~ dgamma(0.001,0.001)
}

```

Figure 5: JAGS model statement for model S.

- 618 to  $\loglambda[i] \sim dt(E.\loglambda[i], \tau.noise, nu)$  and second, we add a prior  $nu \sim$   
 619  $dgamma(2, 0.10)$  for the degrees of freedom parameter (Juárez and Steel, 2010).
- 620 In models S and SH, we assume that year effects  $\gamma_{s,y}$  are conditionally independent and  
 621 normally distributed with precision  $\tau_s^\gamma$ ; the precision is allowed to vary among strata. The

```

# Year effects #
for(s in 1:nstrata){
  tau.gamma.tiny[s] <- tau.gamma[s]*0.0001
  gamma[s,y.baseline] ~ dnorm(stratum[s],tau.gamma.tiny[s])
  for(y in 1:(y.baseline-1))
    gamma[s,y.baseline-y] ~ dnorm(gamma[s,y.baseline-y+1],tau.gamma[s])
  for(y in (y.baseline+1):nyears){
    gamma[s,y] ~ dnorm(gamma[s,y-1],tau.gamma[s])
  }
  stratum[s] ~ dnorm(mu.strata,tau.strata)
  tau.gamma[s] ~ dgamma(0.001,0.001)
}
mu.strata ~ dnorm(0.0,1.0E-6)
tau.strata ~ dgamma(0.001,0.001)

```

Figure 6: Changes to slope models (S and SH) to produce difference models (D and DH).

622 expected value of  $\gamma_{s,y}$  is

$$E(\gamma_{s,y}) = S_s + \beta_s (y - y_0) ;$$

623 the intercept  $S_s$  is a baseline abundance parameter for stratum  $s$ ,  $\beta_s$  is a trend parameter,  
 624 and  $y_0$  is a baseline year to center the regression. Baseline abundance and trend parameters  
 625 are allowed to vary by stratum;  $S_s$  and  $\beta_s$  are modeled as random effects, across strata.

626 In models D and DH, we replace the assumption that year effects  $\gamma_{s,y}$  are conditionally  
 627 independent, with the assumption that *differences* in year effects are conditionally independent.  
 628 Thus,  $\gamma_{s,y}$  is normally distributed with mean  $\gamma_{s,y-1}$  and precision  $\tau_s^\gamma$ . As with the slope  
 629 models, we fix a baseline year  $y_0$ , and fix  $E(\gamma_{s,y_0}) = S_s$ , at the stratum mean. Changes needed  
 630 for the JAGS code are illustrated in Figure 6.

631 In all four models, prior distributions for means are assigned flat normal priors, and prior  
 632 distributions for precisions are assigned flat gamma priors.

633 **Appendix S2: Abundance indices for models with t-distributed overdispersion**

634 **effects**

635 *Background*

636 Let  $Y_i$  denote a count obtained from the North American Breeding Survey. We suppose that

637  $Y_i \sim Pois(\lambda_i)$ , with

$$\log(\lambda_i) = A_{s(i)} + \beta_{s(i)}(y(i) - y^*) + \gamma_{s(i),y(i)} + \omega_{o(i)} + \epsilon_i ;$$

638 here,  $A_s$  and  $\beta_s$  are parameters for stratum  $s$ ;  $s(i), y(i)$  and  $o(i)$  are the stratum, year, and

639 observer associated with count  $i$ ; and  $\omega_o$  and  $\epsilon_i$  are random effects.

640 As an index to abundance for stratum  $s$  in year  $y$ , we use the expected value of  $Y_i$  given

641  $s(i) = s$  and  $y(i) = y$ . Without specific assumptions on the distributions of  $\omega_o$  and  $\epsilon_i$ , the

642 double expectation theorem yields

$$E(Y_i | s(i) = s, y(i) = y) = \exp(A_s + \beta_s(y - y^*) + \gamma_{sy}) \times E(\exp(\omega_{o(i)} + \epsilon_i)) , \quad (8)$$

643 provided that  $E(\exp(\omega_{o(i)} + \epsilon_i))$  exists. Assume that  $\omega_o$  and  $\epsilon_i$  are mean zero normal random

644 variables with variances  $\sigma_\omega^2$  and  $\sigma_\epsilon^2$ . Then  $E(\exp(\omega_{o(i)} + \epsilon_i)) = \exp(\frac{1}{2}\sigma_\omega^2 + \frac{1}{2}\sigma_\epsilon^2)$ . The right-hand

645 side of equation (8) is

$$\theta_{sy} = \exp(A_s + \beta_s(y - y^*) + \gamma_{sy} + \frac{1}{2}\sigma_\omega^2 + \frac{1}{2}\sigma_\epsilon^2) , \quad (9)$$

646 which can be calculated as a derived parameter and monitored in fitting the model via

647 MCMC. Parameter  $\theta_{sy}$  is the expected value of a new observation  $Y_i^*$ , a Rao-Blackwellized  
648 version of a posterior prediction (Casella and Robert 1996), and is thus an appropriate index  
649 to abundance. It should be noted that  $\theta_{sy}$  is a parameter, rather than a statistic; point and  
650 interval estimates for the index are based on its posterior distribution.

651 Suppose now that  $\epsilon_i$  has a t-distribution with  $\nu$  degrees of freedom. Remarkably,  $E(\exp(\epsilon_i)) =$   
652  $\infty$ , for any finite value of  $\nu$ , no matter how large. Thus an expected count (like equation (8))  
653 is inappropriate as a representation of typical counts. How then are we to define an index to  
654 abundance for BBS counts, if we choose to model  $\epsilon_i$  as having a t-distribution?

655 We seek an index that is comparable to the expected count index used under assumptions of  
656 normality. Our solution is to regard the right-hand side of equation (9) as an exponentiated  
657 percentile of the normal distribution. Such a definition for the index  $\theta_{sy}$  can then be  
658 generalized to t-distributions, as we show in these notes.

659 Before doing that, it is worth examining alternative indices to abundance. One would be to  
660 simply use the median value of  $\lambda_i$  given  $s(i) = s$  and  $y(i) = y$ . The index could be called the  
661 “median expected count”; it would only involve the first term on the right-hand side of (8).  
662 However, the median expected count can be substantially smaller than the expected count.  
663 It seems more appropriate to exponentiate a percentile other than the median, and to be  
664 guided by the existing indices obtained under assumptions of normality.

665 *Matching T and Normal distributions*

666 If  $\log(\lambda) \sim N(\mu, \sigma^2)$ , then  $E(\lambda) = \exp(\mu + \frac{1}{2}\sigma^2)$ . Writing  $X = \log(\lambda)$ , we can think of this  
667 quantity as  $E(e^X)$  or as  $\exp(X^*)$ , where  $X^* = \mu + \frac{1}{2}\sigma^2$ . Furthermore, we can think of  $\mu$  and  
668  $\sigma$  as the mean and variance, or as the median and the scaling parameter, relative to the

669 standard normal.

670 T-distributions with  $\nu$  degrees of freedom have density  $\pi_{\nu,k}(t) \propto (1 + ((t - m)/k)^2/\nu)^{-(\nu+1)/2}$

671 ; here  $k$  is a scaling parameter, and  $m$  is the median of the distribution. We have noted that

672 if  $X$  has a t-distribution with  $\nu$  degrees of freedom,  $E(\exp(X))$  does not exist. However, if

673  $\nu > 2$ ,  $E(X)$  and  $\text{Var}(X)$  exist:  $E(X) = m$  and

$$\text{Var}(X) = \frac{k^2 \nu}{\nu - 2} , \quad (10)$$

674 where  $k > 0$  is a scale parameter. We could calculate a value  $X^* = \text{mean} + \frac{1}{2}\text{variance}$ , and

675 define our index of abundance by  $\exp(X^*)$ , even though  $E(\exp(X))$  is not defined. However,

676 the non-existence of a variance for  $\nu \leq 2$ , and the potential instability of (10) for  $\nu$  only

677 slightly larger than 2, argue against this approach. Indeed, if  $\nu \leq 1$ , the t-distribution has

678 neither variance nor mean.

679 Instead, we will think of  $X^*$  as the median value of the t-distribution plus one half of a

680 squared scaling parameter, with scaling relative to the standard normal distribution. We will

681 accomplish the scaling by choosing the value  $k = k(\nu)$  that leads to the closest approximation

682 to a standard normal distribution by a central t-distribution with  $\nu$  degrees of freedom

683 and scale parameter  $k$ . The best approximating value of  $k(\nu)$  will be that which minimizes

684 Kullback-Leibler (KL) divergence from the standard normal distribution.

685 For each of 250 values of  $\nu$ , evenly spaced on the interval [0.5, 40], we minimized the KL

686 divergence

$$\int_{-\infty}^{\infty} \log \left( \frac{\phi(t)}{\pi_{\nu,k}(t)} \right) \phi(t) dt , \quad (11)$$

687 as a function of  $k$ ; here  $\phi(\cdot)$  is the density of a standard normal distribution. We approximated  
688 the integral in (11) by a sum on  $10^6$  grid points over  $[-10,10]$ .

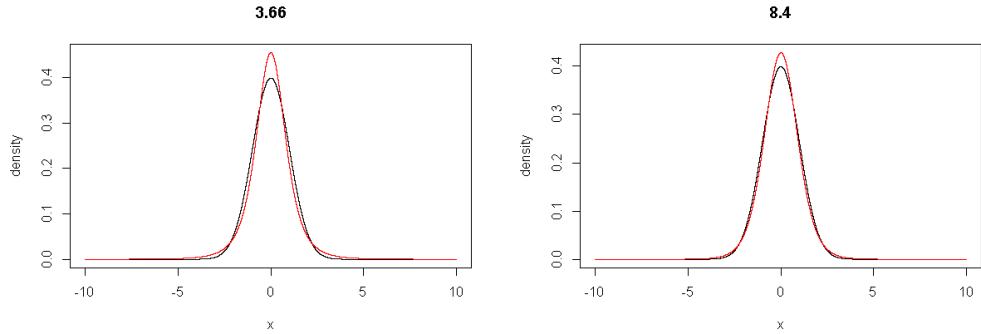


Figure 7: Standard normal (black) and best approximating  $t_{\nu,k}$  distributions (red), for  $\nu = 3.66, 8.40$ .

689 Two examples of the best matching densities are given in Figure (7). The relation between  
690 the best scale parameter and the degrees of freedom is very close to linear on the log-logistic  
691 scale ( $R^2 = 0.9989$ ); the optimal values of  $k$  and fitted regression are given in Figure (8). The  
692 regression relation is

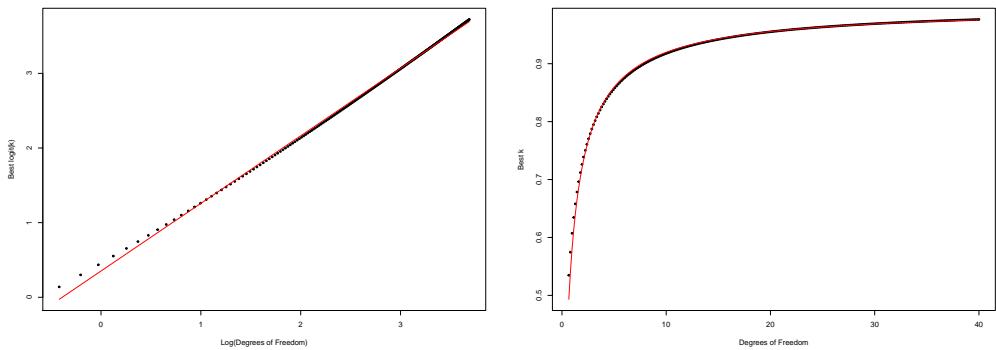


Figure 8: Fitted relation (red curves) between best  $k$  (black dots) and degrees of freedom.

$$k(\nu) = \frac{1.422 \nu^{0.906}}{1 + 1.422 \nu^{0.906}} . \quad (12)$$

693 The value  $k(\nu)$  is the scale parameter that most closely matches a  $t_{\nu,k}$  distribution to a  
 694 standard normal distribution, in terms of K-L divergence (equation (11)). From (10) it is  
 695 anticipated that  $k(\nu) \approx \sqrt{\frac{\nu-2}{\nu}}$ ; indeed, the ratio of the two is 1.04889 for  $\nu = 7.5$ , decreasing  
 696 to a minimum of 0.9992 at  $\nu = 151$ , then tending to 1 as  $\nu \rightarrow \infty$ .

697 In JAGS code, we write  $X \sim dt(m, tau, nu)$ . The parameter `tau` is a squared reciprocal  
 698 of the basic scaling parameter, and must be adjusted by  $k(\nu)$ : the desired scaling value is  
 699  $(1/\sqrt{tau})/k(nu)$ .

### 700 Appendix S3: The need for and usefulness of leave-one-out calculations

701 Suppose data are a sample of  $n = 20$  independent Poisson variables with expected value  $\lambda$ .  
 702 Our data set happens to have an observation  $Y_i = 11$ , but the sample mean is a good bit  
 703 lower, say  $\bar{Y} = 5.3$ . We want to calculate  $\Pr(Y = 11|\lambda)$  to see whether the observation of 11  
 704 is consistent with the Poisson model. The problem is, we don't know  $\lambda$ , and the fact that one  
 705 of our observations was  $Y_i = 11$  suggests that the sample mean might be a bit high, with the  
 706 result that the calculated value  $\Pr(Y = 11|\lambda = 5.3) = 0.0116$  is too large. Omitting the value  
 707  $Y_i = 11$  from the data set, the sample mean calculated from the remaining 19 observations is  
 708  $(20 \times 5.3 - 11)/19 = 5.0$ , and  $\Pr(Y = 11|\lambda = 5.0) = 0.0082$ . The original value, based on the  
 709 full data set, was 40% larger. Not surprisingly, the inclusion of  $Y_i$  in the estimation of  $\Pr(Y_i)$   
 710 makes the observation look more consistent with the model.

711 The very nature of maximum likelihood estimation guarantees such behavior, in evaluating

712 the fit of a model to a data set. The likelihood for  $\theta^M$  is the product  $\prod_i \Pr(Y_i|\theta^M, M)$ ,  
713 considered as a function of  $\theta^M$ . We would like to evaluate this function at the true value of  
714  $\theta^M$ . By definition the MLE provides the largest possible value of the product, a value which  
715 is of necessity greater than or equal to what we'd get using the true value.

716 The same effect occurs in a Bayesian analysis. We continue with the example of estimating  
717 probabilities for a Poisson random variable with mean  $\lambda$ . We will write  $f_k(\lambda) = \Pr(Y =$   
718  $k|\lambda) = \exp(-\lambda) \lambda^k/k!$ , our notation emphasizing that the quantity of interest is a derived  
719 parameter, i.e., a function of  $\lambda$ .

720 Given a sample of size  $n$ , denoted by  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)$ , and a prior distribution for  $\lambda$ ,  
721 inference will be based on the posterior distribution  $[f_k(\lambda)|\mathbf{Y}]$ , which can be obtained by  
722 transformation of the posterior distribution  $[\lambda|\mathbf{Y}]$ . In particular, we will suppose that a  
723 flat gamma prior has been chosen for  $\lambda$ , so that  $[\lambda|\mathbf{Y}] = \text{Gamma}(\lambda; S, n)$ , where  $S = \sum_i Y_i$ .  
724 Sampling  $\lambda$ 's from this posterior distribution and calculating  $f_k(\lambda)$  produces a sample of  
725  $[f_k(\lambda)|\mathbf{Y}]$ .

726 Though our estimation procedure is Bayesian, we evaluate its frequency properties, i.e., the  
727 typical values of  $[f_k(\lambda)|\mathbf{Y}]$  across replicate samples  $\mathbf{Y}$  with a common (generating) value of  
728  $\lambda$ . Frequentist evaluations of Bayesian procedures were advocated by Rubin (1984); Bayesian  
729 procedures with good frequentist operating properties are described as well-calibrated.

730 In what follows, we present 3 calculations. First, we average  $[f_k(\lambda)|\mathbf{Y}]$  over an unconstrained  
731 set of replicate data sets  $\mathbf{Y}$ . We denote the resulting out-of-sample mean posterior distribution  
732 by  $P_{os}^k(\lambda)$ . Next, we average  $[f_k(\lambda)|\mathbf{Y}]$  over a constrained set of replicate data sets  $\mathbf{Y}$ , chosen  
733 to include an observation  $Y_i = k$ . We denote the resulting within-sample mean posterior  
734 distribution by  $P_{ws}^k(\lambda)$ . Finally, we average a leave-one-out distribution over the same

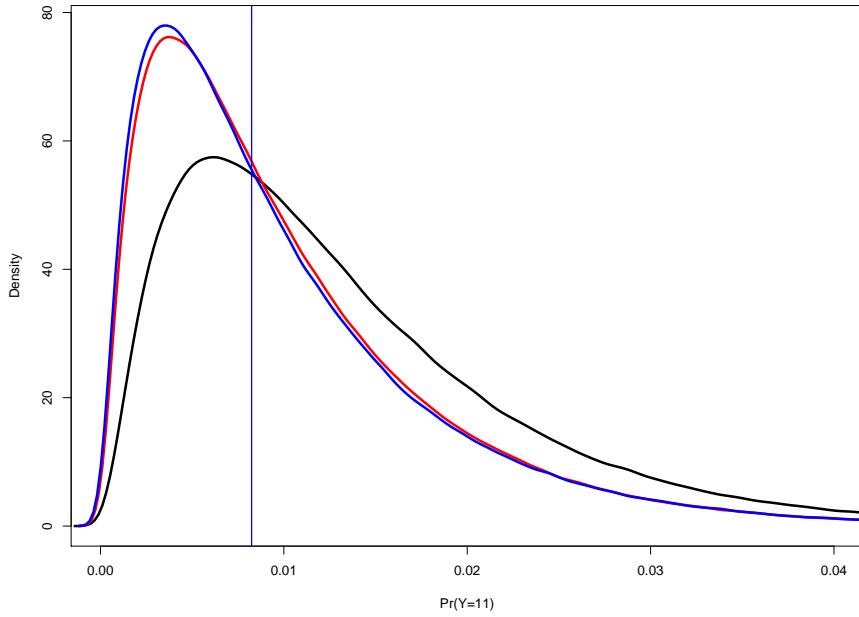


Figure 9: Average values of  $\Pr(Y = 11|\mathbf{Y})$  taken over random samples  $\mathbf{Y}$  of 20 Poisson(5) random variables. Using notation defined in text, red curve is density for  $P_{os}^{11}(\lambda)$ , black curve is for  $P_{ws}^{11}(\lambda)$ , and blue curve is for  $P_{Loo}^{11}(\lambda)$ .

735 constrained set of replicate data; we denote the resulting within-sample mean leave-one-out  
 736 distribution by  $P_{Loo}^k(\lambda)$ .

737 Figure 9 displays the mean value of  $[f_k(\lambda)|\mathbf{Y}]$  for  $k = 11$ , averaged over  $10^6$  sampled data  
 738 sets with  $n = 20$  and  $\lambda = 5$  (red curve). This density can be thought of as the typical basis  
 739 of inference about  $\Pr(Y = 11|\lambda)$  for a sample of size 20, given  $\lambda = 5$ . We refer to this as  
 740 an out-of-sample mean, because (in contrast to subsequent calculations) the sampling of  
 741 replicate data sets was done without reference to any outcomes in a given data set. Inference  
 742 based on  $P_{os}^k(\lambda)$  appears reasonably well calibrated, with 51.4% of the mass in the density  
 743 lying to the left of the true value 0.00824 (indicated by vertical blue line).

744 Now suppose that the analyst's interest in  $\Pr(Y = 11|\lambda)$  is sparked by the occurrence of an  
 745 observation  $Y_i = 11$  in his dataset. If  $\lambda = 5$ , an observation  $Y_i = 11$  stands out as large in

746 a sample of size 20. The sample itself is atypical; its sample mean is likely larger than the  
747 true mean, and estimates of  $\Pr(Y = 11|\lambda)$  will be too large. We need to evaluate the average  
748 value of  $[f_{11}(\lambda)|\mathbf{Y}]$  over replicate data sets including an observation equal to 11. So once  
749 again we generated  $10^6$  sampled data sets with  $n = 20$  and  $\lambda = 5$ ; this time, however, we used  
750 rejection sampling to require that at least one observation in the data set was equal to 11.

751 The result is the within-sample mean posterior distribution  $P_{ws}^k(\lambda)$ , plotted in black in Figure  
752 9. Only 35.8% of the mass of  $P_{ws}^k(\lambda)$  lies to the left of the true value: the estimates tend to  
753 be too large, as anticipated. The mean value is 30.9% larger than the mean for  $P_{ws}^k(\lambda)$ . Using  
754 sampled values to compute their own probabilities leads to poorly-calibrated inferences.

755 The solution is to omit the value  $Y_i = 11$  from the data set and treat what's left (denoted  
756  $\mathbf{Y}_{-i} = \mathbf{Y} \setminus Y_i$ ) as though it were a random sample of size  $n - 1$ . We thus sample  $\lambda \sim$   
757  $\text{Gamma}(\lambda; S - Y_i, n - 1)$  and compute  $f_{11}(\lambda)$ . With a slight abuse of notation, we describe  
758 the collection of sampled values as a sample of  $[f_{11}(\lambda)|\mathbf{Y}_{-i}]$ . To evaluate the calibration of  
759 this procedure, we once again use rejection sampling to generate  $10^6$  data sets with at least  
760 one observation equal to 11. The average value of  $[f_{11}(\lambda)|\mathbf{Y}_{-i}]$  for these data sets is  $P_{Loo}^k(\lambda)$ ,  
761 plotted in blue in Figure 9; it closely approximates the density for  $P_{os}^k(\lambda)$ .

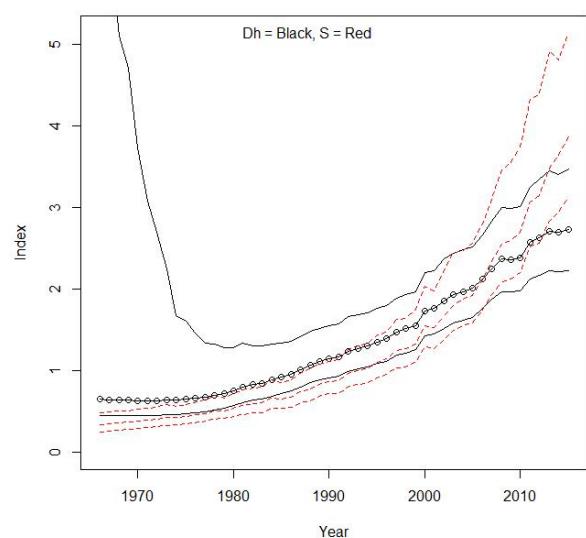
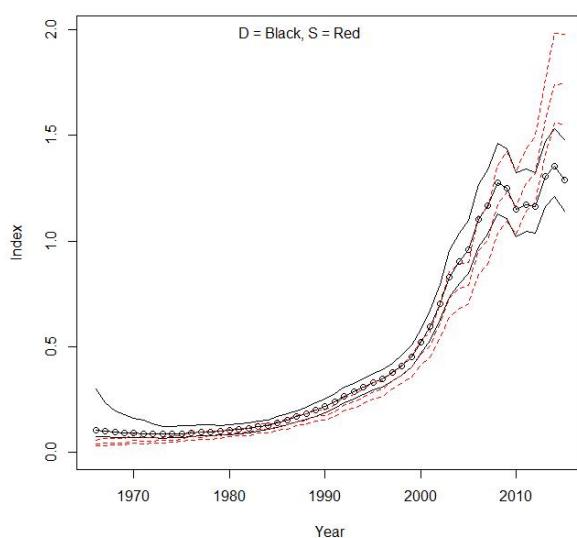
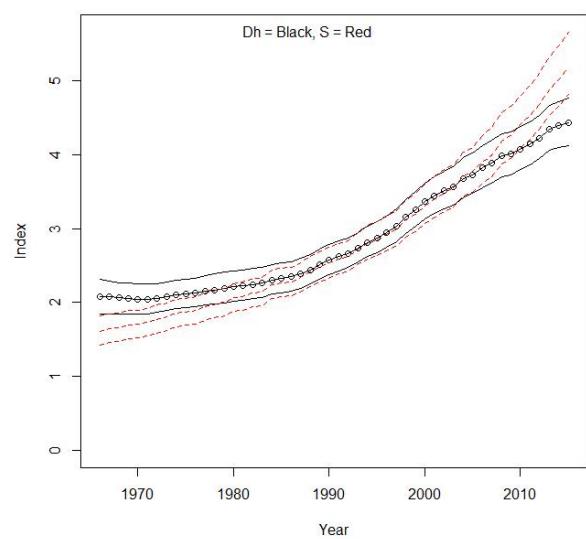
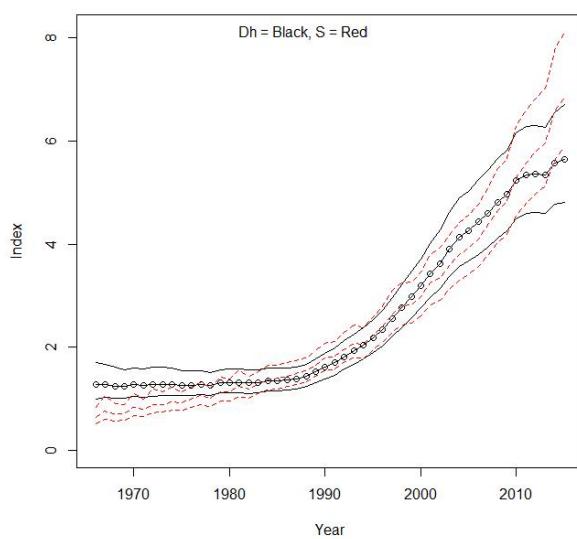
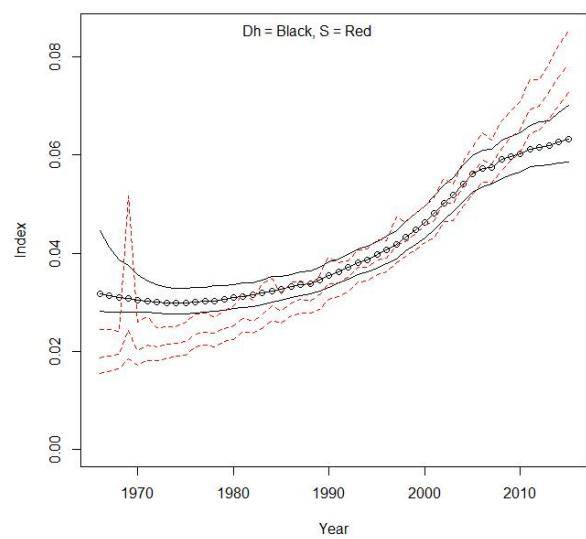
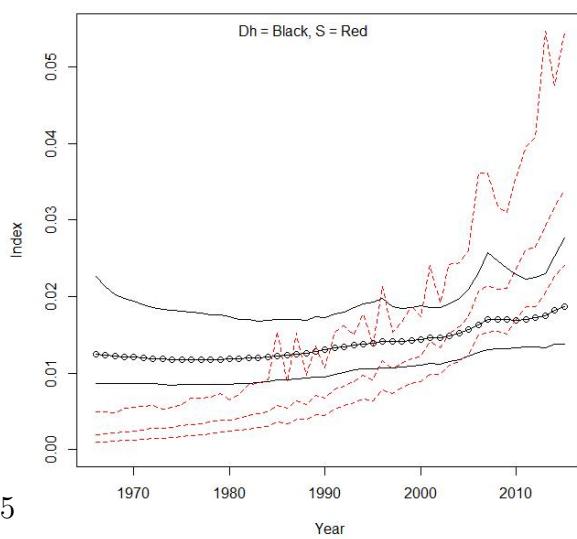
762 While use of the posterior distribution  $[f_k(\lambda)|\mathbf{Y}]$  has good out-of-sample properties, its within-  
763 sample tendency is to overestimate. Using an observation to predict the probability of one just  
764 like it isn't a good idea. The solution to this problem is to use leave-one-out calculations. Even  
765 when evaluated using within-sample simulations, the leave-one-out procedure has frequency  
766 properties approximating the out-of-sample properties of the full posterior inference.

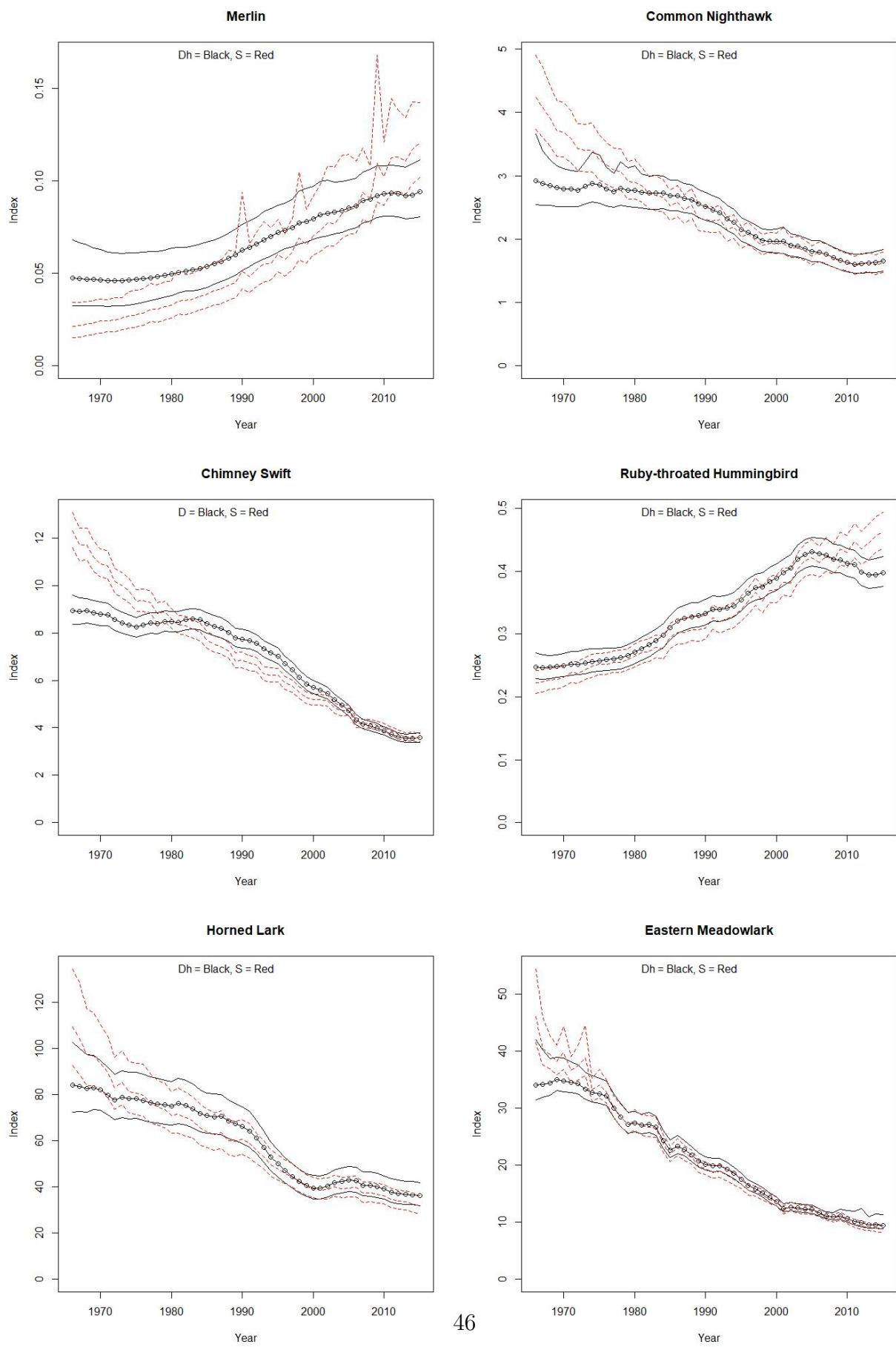
767 **Appendix S4: Results for selected vs traditional model**

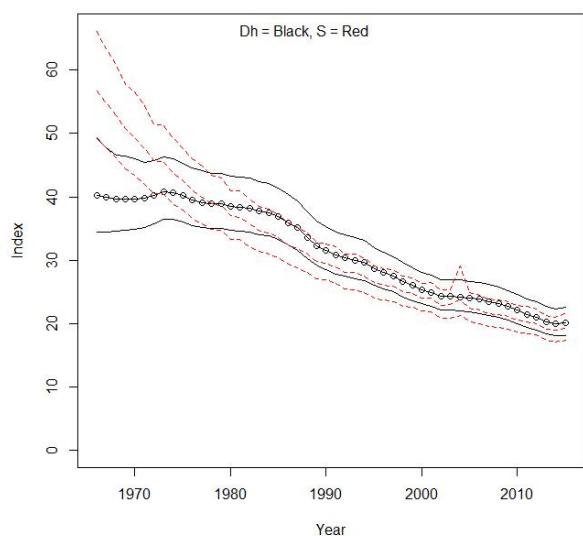
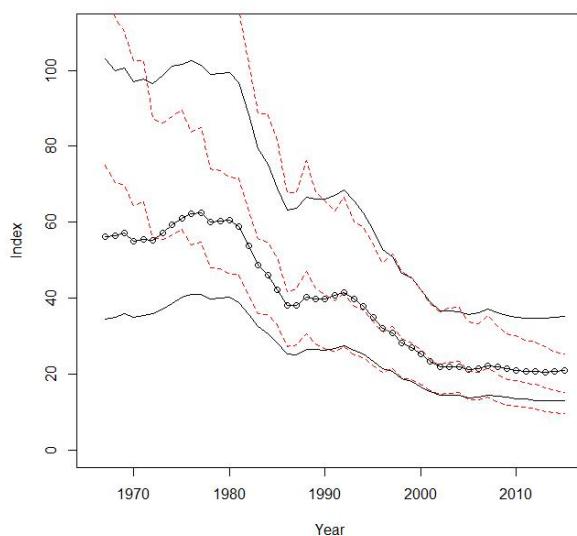
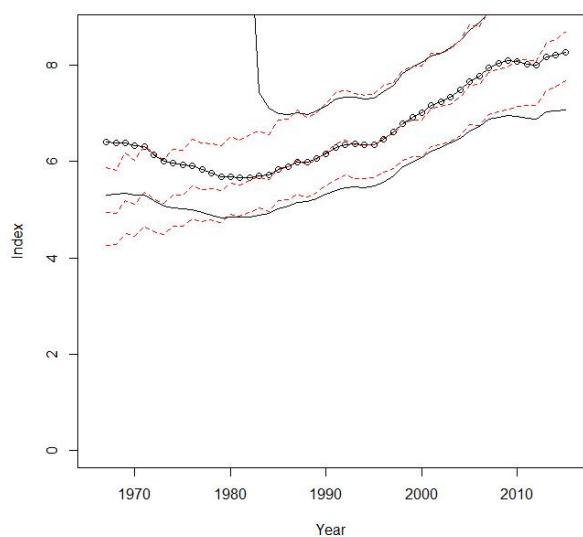
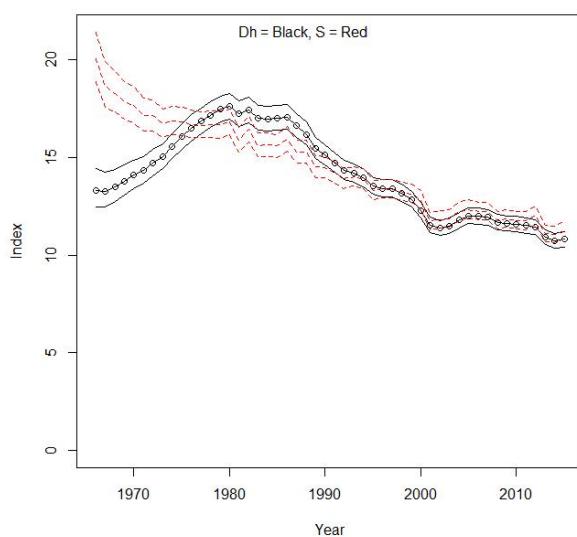
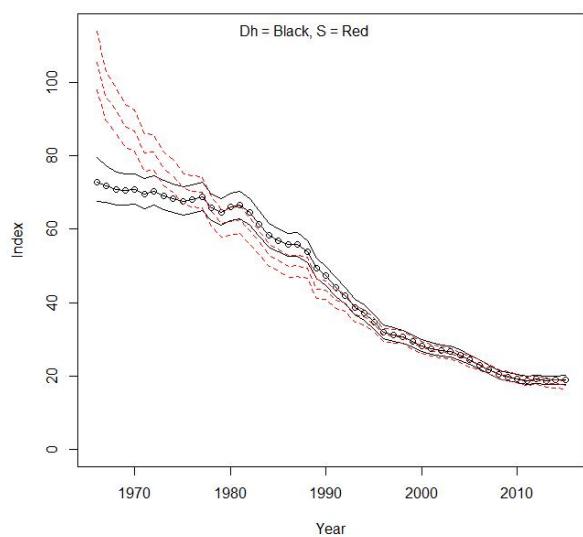
768 Excel spreadsheet Table of Trend Results 1970–2015.xlsx summarizes population trend  
769 (%/yr) and 95% credible interval limits (2.5<sup>th</sup> and 97.5<sup>th</sup> percentiles) estimated at the survey-  
770 wide scale from BBS data for 548 species. We present results for the default model (S) and  
771 the selected model. We also present numbers of survey routes used in the analysis, and  
772 Relative Abundances (Rel. Abun., defined as the estimated annual index (birds/route) for  
773 the species in the year 1990) based on the default and selected model. Core species results are  
774 based on the interval 1970–2015; non-core species results are based on the interval 1993–2015.

775 **Appendix S5: Annual indices for 18 (of 426) core species with discrepant results  
776 between selected model and traditional model**

777 Annual indices for 18 species for which trend results differed between the dafault model (S)  
778 and the selected model. Annual indices and 95% credible intervals for the default model  
779 are portrayed in red stippled; indices for selected model are portrayed in solid black. For 3  
780 species (Wild Turkey (*Meleagris gallopavo*), Chimney Swift (*Chaetura pelagica*), and Eurasian  
781 Collared-Dove (*Streptopelia decaocto*)) model D was selected; model DH was selected for  
782 the remaining 15 species (Sandhill Crane (*Antigone canadensis*), Turkey Vulture (*Cathartes*  
783 *aura*), Black Vulture (*Coragyps atratus*), Cooper's Hawk (*Accipiter cooperii*), Merlin (*Falco*  
784 *columbarius*), Peregrine Falcon (*Falco peregrinus*), Common Nighthawk (*Chordeiles minor*),  
785 Ruby-throated Hummingbird (*Archilochus colubris*), Horned Lark (*Eremophila alpestris*),  
786 Eastern Meadowlark (*Sturnella magna*), Brewer's Blackbird (*Euphagus cyanocephalus*),  
787 Chestnut-collared Longspur (*Calcarius ornatus*), Western Tanager (*Piranga ludoviciana*),  
788 Barn Swallow (*Hirundo rustica*), and House Sparrow (*Passer domesticus*).)

**Sandhill Crane****Wild Turkey****Turkey Vulture****Black Vulture****Cooper's Hawk****Peregrine Falcon**



**Brewer's Blackbird****Chestnut-collared Longspur****Western Tanager****Barn Swallow****House Sparrow****Eurasian Collared-Dove**