

1 **Comparing the performance of three data weighting methods when**
2 **allowing for time-varying selectivity**

3
4 Haikun Xu^{1*,2} (hkxu@iattc.org), James T. Thorson^{3,4} (James.Thorson@noaa.gov), Richard D.
5 Methot⁵ (Richard.Methot@noaa.gov)

6
7 ¹ *Inter-American Tropical Tuna Commission, 8901 La Jolla Shores Drive, La Jolla, CA 92037,*
8 *USA*

9 ² *Previously at School of Aquatic and Fishery Sciences, University of Washington, Box 355020,*
10 *Seattle, WA 98105, USA*

11 ³ *Habitat and Ecosystem Process Research program, Alaska Fisheries Science Center, National*
12 *Marine Fisheries Service, National Oceanic and Atmospheric Administration, 7600 Sand Point*
13 *Way NE, Seattle, WA 98115*

14 ⁴ *Previously at Fishery Resource Analysis and Monitoring Division, Northwest Fisheries Science*
15 *Center, National Marine Fisheries Service, National Oceanic and Atmospheric Administration,*
16 *2725 Montlake Blvd. East, Seattle, WA 98112, USA*

17 ⁵ *NOAA Senior Scientist for Stock Assessments, National Marine Fisheries Service, National*
18 *Oceanic and Atmospheric Administration, 2725 Montlake Blvd. East, Seattle, WA 98112, USA*

19
20
21
22
23
24
25
26
27
28
29 ** Corresponding author address:*

30 Haikun Xu, Inter-American Tropical Tuna Commission,

31 8901 La Jolla Shores Drive, La Jolla, CA 92037, USA

32 Telephone: +1 8583342824; E-Mail: hkxu@iattc.org

33 **Abstract**

34 How to properly weight composition data is an important ongoing research topic for fisheries stock
35 assessments and multiple methods for weighting composition data have been developed. Although
36 several studies indicated that properly accounting for time-varying selectivity can reduce
37 estimation biases in population biomass and management-related quantities, no study to date has
38 compared the performance of widely-used data weighting methods when allowing for time-
39 varying selectivity. Here, we conducted four simulation experiments for this topic, aiming to
40 provide guidance on weighting age-composition data given time-varying selectivity. The first
41 simulation experiment showed that over-weighting should be avoided in general and even when
42 estimating time-varying selectivity. The second simulation experiment compared three data
43 weighting methods (McAllister-Ianelli, Francis, and Dirichlet-multinomial), within which the
44 Dirichlet-multinomial method outperformed the other two methods when selectivity is time-
45 varying. The third and fourth simulation experiments further showed that given time-varying
46 selectivity, the Dirichlet-multinomial method still performed well when age-composition data
47 were over-dispersed and when the level of selectivity variation needed to be simultaneously
48 estimated. Our simulation results support using the Dirichlet-multinomial method when estimating
49 time-varying fishery selectivity. Also, the simulation results suggest that improving stock
50 assessments by accounting for time-varying selectivity requires simultaneously addressing data
51 weighting and time-varying selectivity.

52 Keywords: Data weighting; time-varying selectivity; Dirichlet-multinomial method; age-
53 composition data

54 **1. Introduction**

55 Fisheries managers use stock assessment models to predict the likely impact of alternative
56 management actions on fishery sustainability. In many jurisdictions worldwide, fisheries managers
57 are recommended or required to manage fishery catches and population abundance in accordance
58 with management targets or limits that are determined from stock assessment models (Methot
59 2009). Accurate predictions of likely management impacts require stock assessment models to
60 appropriately approximate biological processes including growth, mortality, maturity, and
61 reproduction. To estimate these different processes, modern assessment models typically fit to a
62 wide range of data sources including abundance indices, subsamples of age/length/sex-
63 composition in fishery-independent surveys or fishery operations, and total fishery landings
64 (Maunder and Punt 2013).

65 Composition data from a fishery are usually not independent between ages and contain a
66 reduced amount of information than they would do if sampled independently (Francis 2011,
67 Maunder 2011, Thorson 2014). Due to, for example, age- or length-specific behaviors such as
68 schooling and aggregating, the age and length of fish from the same set are more similar than from
69 different sets. Namely, composition samples are positively correlated among adjacent age or length
70 bins, contradicting the assumption of random sampling in the widely used multinomial distribution
71 for composition data (Francis 2011). This phenomenon, which is referred to as “over-dispersion”,
72 increases the variance in composition samples and decreases the effective sample size. The
73 weighting of composition data in stock assessment models is positively related to the effective
74 sample size, which is used by stock assessment scientists to accommodate unknown observation
75 error and model mis-specification.

76 Stock assessment models will often estimate different values for stock status and productivity
77 when fitted to different subsets of available data (Maunder et al. 2017). In particular, inferring
78 trends in population abundance from age- and length-composition sampling depends upon correct
79 specification of many biological processes including mortality, growth, and availability to survey
80 or fishery operations, and mis-specification of these processes will cause information in age- and
81 length-composition data to be biased with respect to true trends in abundance (Minte-Vera et al.
82 2017). For this reason and others, several papers have suggested that age- and length-composition
83 data should be “down-weighted” relative to abundance index data whenever the two provide
84 conflicting information about abundance trends (Francis 2011, 2014). Widely-used methods for
85 weighting composition data include the methods by McAllister and Ianelli (1997), Francis (2011),
86 and the linear parameterization of the Dirichlet-multinomial (D-M) likelihood (Thorson et al.
87 2017). These and other methods all have in common that they down-weight age- and length-
88 composition data more when the assessment model predictions and available data are greatly
89 different, and down-weight less (or even up-weight) when predictions and data match well.
90 However, these methods also differ in well-documented respects: the McAllister-Ianelli (M-I) and
91 Francis methods require iteratively fitting a stock assessment model and do not characterize model
92 uncertainty caused by estimating data-weights, while the D-M method can be efficiently estimated
93 as a model parameter with associated measure of uncertainty (Francis 2017, Thorson 2018).

94 One main reason for down-weighting composition data is that stock assessment models explain
95 process error as observation or sampling error (Maunder and Piner 2017). For example, when a
96 time-varying selectivity is mis-specified to be time-invariant, the stock assessment model
97 attributes the discrepancy between observed and expected compositions solely to the error in the
98 composition sampling process. As an alternative, analysts may instead revise a stock assessment

99 model such that it is better able to predict available data. There are many biological processes that
100 could cause the proportion of an age/length/sex-composition that are selected by a given fishery
101 or survey operation to vary over time, including spatial patterns in fishery effort (Sampson and
102 Scott 2012), environmentally-driven changes in vertical distribution (Kotwicky et al. 2015), or
103 spatial redistribution among well- and poorly-sampled habitats (Thorson et al. 2013a). In general,
104 these processes will cause “model mis-specification”, wherein a model cannot match available
105 data even if unlimited or perfect data are available. In these cases, a stock assessment can estimate
106 additional fixed or random effects representing time-varying selectivity, and this will generally
107 increase the match between available data and model predictions (Lowe et al. 2017, Martell and
108 Stewart 2014, Xu et al. 2018). In fact, a number of simulation studies have shown that properly
109 accounting for time-varying selectivity can reduce estimation biases in population biomass and
110 management-related quantities (Stewart and Martell 2014, Stewart and Monnahan 2017, Thorson
111 and Taylor 2014, Xu et al. 2018).

112 By increasing the match between model predictions and data, estimating time-varying
113 selectivity will clearly impact the degree of data weighting estimated by different methods.
114 However, no study to date has compared the performance of widely-used data weighting methods
115 in the case where the assessment model estimates time-varying selectivity. Under the assumption
116 of constant selectivity, Maunder (2011) showed that estimating the effective sample size of
117 composition data led to an improvement over using the nominal sample size (the number of fish
118 sampled each year) if the corresponding true selectivity varied from year to year. Under the
119 assumption of time-varying selectivity, in comparison, estimating the effective sample size of
120 composition data (namely, weighting composition data) is more problematic because the true level
121 of variation in selectivity is unknown. Using a simulation approach, Stewart and Monnahan (2017)

122 explored the effects of data weighing on the performance of models with or without process error
123 in selectivity. Based on simulation results, they concluded that assessment models should allow
124 for process error in selectivity and should not excessively down-weight composition data.

125 The main objective of this paper was to compare the performance of three data weighting
126 methods when allowing for time-varying selectivity. Previous studies (e.g., Hulson et al. 2012) has
127 evaluated the performance of several data weighting methods using simulations, but our study is
128 the first to use simulation to compare the performance of data weighting methods in assessment
129 models that estimate time-varying selectivity. We first conducted a simulation experiment to
130 evaluate the sensitivity of model performance to the extent to which fisheries age-composition
131 data are weighted in assessment models both with and without process error in selectivity, given
132 that the true selectivity is time-varying. This experiment aimed to answer the question: what are
133 the consequences of under- or over-weighting age-composition data when process error in
134 selectivity is ignored or estimated? We then conducted three simulation experiments to compare
135 the performance of three (M-I, Francis, and D-M) data weighting methods given that the true
136 selectivity is time-varying, aiming to address the following questions:

- 137 1) Which data weighting method performs best when the assessment model estimates time-
138 varying selectivity?
- 139 2) How is the performance of the best data weighting method degraded owing to the over-
140 dispersion in age-composition data?
- 141 3) Can we simultaneously weight age-composition data and estimate the selectivity variation
142 penalty in stock assessments?

143

144 **2. Materials and methods**

145 In this paper, we compared three methods for weighting age-composition data based on
146 simulation experiments that were undertaken by modifying an age-structured simulation-
147 estimation package *CCSRA* (Thorson and Cope 2015). We first described the basic structure and
148 hypotheses for the operating model (OM), sampling model (SM), and estimation model (EM) used
149 in our simulation experiments. We then described in detail how the OM, SM, and EM were
150 configured and how model performance was evaluated in each simulation experiment. In each
151 simulation experiment, the OM simulated the true population dynamics from which the SM
152 generated observation data. The EM was then fitted to the generated observation data and model
153 performance was evaluated by comparing the estimates of population attributes that the EM
154 provided with the corresponding true values that the OM simulated.

155 *2.1. Simulation models*

156 *2.1.1. Operating model*

157 The OM was an age-structured model (Table 1) that allows fishery selectivity to vary either
158 independently or correlated from a specified parametric functional form. It was used in this study
159 to simulate the true population dynamics for two species, Pacific hake (*Merluccius productus*) and
160 Pacific sardine (*Sardinops sagax*), that correspond to a “periodic” and “opportunistic” type of life
161 history, respectively (Table 2). The level of recruitment variation (σ_R in Eq. T1.2) was specified
162 to be low (0.4) for Pacific hake and to be either low (0.4) or high (0.8) for Pacific sardine, in order
163 to compare the performance of data weighting methods under contrasting levels of recruitment
164 variation. A higher level of recruitment variation caused a larger contrast in each year’s age-
165 composition observation. The two types of life history with differing recruitment assumptions were
166 hereafter referred to as hake-low, sardine-low, and sardine-high. The OM included one fishery and

167 the selectivity of which in age a and year t was specified to be a product of a parametric (logistic)
 168 form and a random deviation term away from the parametric form:

$$169 \quad S_{a,t} = \frac{1}{1 + e^{-S_{slope}(a-S_{50})}} \times e^{\varepsilon_{a,t}} \quad (1)$$

170 Particularly, the non-parametric deviation term ($\varepsilon_{a,t}$), which can be treated as a process error in
 171 fishery selectivity, was specified to follow a two-dimensional AR(1) process:

$$172 \quad \text{vec}(\boldsymbol{\varepsilon}) \sim \text{MVN}(\mathbf{0}, \sigma_S^2 \mathbf{R} \otimes \tilde{\mathbf{R}}) \quad (2a)$$

$$173 \quad \mathbf{R}_{a,\bar{a}} = \rho_a^{|a-\bar{a}|} \quad (2b)$$

$$174 \quad \tilde{\mathbf{R}}_{t,\bar{t}} = \rho_t^{|t-\bar{t}|} \quad (2c)$$

175 where $\sigma_S (>0)$ is the standard deviation of selectivity deviations that controls the degree of
 176 variation in fishery selectivity and $\rho_a (-1 < \rho_a < 1)$ and $\rho_t (-1 < \rho_t < 1)$ are two AR(1) coefficients that
 177 control the degree to which selectivity deviations are autocorrelated in age and time, respectively.
 178 The deviations of fishery selectivity are identical and independent (IID) when ρ_a and ρ_t are both
 179 zeroes because this specification simplifies Eq. 2a to be

$$180 \quad \varepsilon_{a,t} \sim N(0, \sigma_S^2) \quad (3)$$

181 We explored four OMs with differing autocorrelation cases under a moderate level of
 182 selectivity variation:

- 183 1. OM1 (“*Independent*”). The deviations of fishery selectivity are independent ($\rho_a = 0; \rho_t =$
 184 $0; \sigma_S = 0.4$);
- 185 2. OM2 (“*Time-correlated*”). The deviations of fishery selectivity are highly autocorrelated in
 186 time ($\rho_a = 0; \rho_t = 0.8; \sigma_S = 0.4$);

- 187 3. OM3 (“Age-correlated”). The deviations of fishery selectivity are highly autocorrelated in age
188 ($\rho_a = 0.8; \rho_t = 0; \sigma_S = 0.4$);
- 189 4. OM4 (“Age- and time-correlated”). The deviations of fishery selectivity are highly
190 autocorrelated in both age and time ($\rho_a = 0.8; \rho_t = 0.8; \sigma_S = 0.4$).

191 We used the *mvrnorm* function in the MASS R package (version 7.3-50, Venables and Ripley
192 2002) to simulate the autocorrelated process error in fishery selectivity. Estimating selectivity
193 deviations is usually difficult for the youngest and oldest age groups due to a lack of adequate age-
194 composition samples for those age groups, so in the OM we assumed that $\varepsilon_{a,t} = \varepsilon_{2,t}$ for $a < 2$ and
195 $\varepsilon_{a,t} = \varepsilon_{7,t}$ for $a > 7$, namely, $\text{vec}(\boldsymbol{\varepsilon}) = (\varepsilon_{2,1}, \dots, \varepsilon_{2,T}, \varepsilon_{3,1}, \dots, \varepsilon_{3,T}, \dots, \varepsilon_{7,1}, \dots, \varepsilon_{7,T})'$. Due to this
196 assumption, the simulated time-varying selectivity cannot be dome-shaped (for $a > 7$, $\varepsilon_{a,T} \equiv \varepsilon_{7,T}$).
197 The parametric selectivity profile as well as the associated variability (induced by the random
198 deviation term) for Pacific hake and Pacific sardine were compared in Figure 1. For both species,
199 we set the plus-group (A) and last simulation year (T) to be 15 and 20, respectively. Fishing
200 mortality was simulated according to an effort-dynamics model (T1.6; more details in Thorson et
201 al. (2013b)) that was used to generate contrast in spawning biomass (SB): the fishery (Fig. 2, left
202 column) drove SB down to about 40% (see Table 2 for parameter values) of the unfished level
203 over 20 years (Fig. 2, right column). A detailed description of how the life history parameters were
204 derived for the two types of life histories (Pacific hake and Pacific sardine) also can be found in
205 Thorson and Cope (2015).

206 2.1.2. Sampling model

207 The sampling model generated the following observation data from the true population
208 dynamics specified by the OM:

- 209 • Fishery total catch in weight, which was assumed to be known without error.
- 210 • Fishery index of abundance (I), which was assumed to be log-normally distributed with a
- 211 coefficient of variation of CV_{abund} and catchability of $q : \ln(I_t) \sim N(\log(qB_t), \ln(1 +$
- 212 $CV_{abund}^2))$, where $B_t = \sum_{a=0}^A N_{a,t} w_a S_{a,t}$ is the exploitable biomass in year t .
- 213 • Fishery age-composition data (\mathbf{A}), which was assumed to be drawn from a multinomial
- 214 distribution with a sample size of $n_{true} : \mathbf{A}_t \sim \text{Multinomial}(\mathbf{P}_t, n_{true})$, where $\mathbf{P}_t = \mathbf{C}_t /$
- 215 $\sum_{a=0}^A C_{a,t}$ is the true age-composition proportion in year t .

216 We assumed that both the index of abundance and age-composition data were informative
 217 ($CV_{abund} = 0.1$ and $n_{true} = 200$) and were available every year during the simulation period.
 218 Therefore, model performance was not limited by low-quality data and should be primarily
 219 determined by how properly the fishery age-composition data (\mathbf{A}_t) were weighted.

220 2.1.3. Estimation model

221 The estimation model had the same population dynamics as the operating model except
 222 whether and how fishery selectivity varied over age and time. Three EMs with differing selectivity
 223 specifications were considered in each simulation experiment:

- 224 • EM1 (“zero deviations”). Selectivity of the fishery was specified to be constant by fixing
- 225 $\hat{\sigma}_S$ as zero: $\hat{S}_a = \frac{1}{1 + e^{-\hat{s}_{slope}(a - \hat{s}_{50})}}$. This specification for fishery selectivity is still a
- 226 common practice in stock assessments.
- 227 • EM2 (“IID deviations”). Selectivity of the fishery was specified to be age- and time-
- 228 varying and the deviations of which were specified to be identical and independent of age
- 229 and time: $\hat{S}_{a,t} = \frac{1}{1 + e^{-\hat{s}_{slope}(a - \hat{s}_{50})}} \times e^{\hat{\epsilon}_{a,t}}$, where $\hat{\epsilon}_{a,t} \sim N(0, \hat{\sigma}_S^2)$.

230 • EM3 (“*AR deviations*”). Selectivity of the fishery was specified to be age- and time-varying
 231 and the deviations of which were specified to be autocorrelated: $\hat{S}_{a,t} = \frac{1}{1+e^{-\hat{S}_{slope}(a-\hat{S}_{50})}} \times$
 232 $e^{\hat{\varepsilon}_{a,t}}$, where $\text{vec}(\hat{\boldsymbol{\varepsilon}}) \sim \text{MVN}(\mathbf{0}, \hat{\sigma}_S^2 \mathbf{R} \otimes \tilde{\mathbf{R}})$, $\mathbf{R}_{a,\bar{a}} = \hat{\rho}_a^{|a-\bar{a}|}$, and $\tilde{\mathbf{R}}_{t,\bar{t}} = \hat{\rho}_t^{|t-\bar{t}|}$.

233 This study was focused on two data weighting issues in stock assessments: sensitivity of model
 234 performance to data weighting and which data weighting method performs better when estimating
 235 age- and time-varying selectivity. In some simulation experiments, we assumed that the hyper-
 236 parameters for selectivity deviations ($\hat{\sigma}_S$, $\hat{\rho}_a$, and $\hat{\rho}_t$) were known without error; in other simulation
 237 experiments, we estimated these hyper-parameters. We included both simulation experiments to
 238 determine model performance either in an idealized case (when these hyper-parameters are known)
 239 or in a more realistic case (when they must be estimated).

240 The age-composition data from the fishery were assumed to be drawn from a multinomial
 241 distribution with an estimated effective sample size of n_{eff} . It specified the extent to which the
 242 fishery age-composition data were weighted:

$$243 \quad \mathbf{A}_t \sim \text{Multinomial}(\hat{\mathbf{P}}_t, n_{eff}) \quad (4)$$

244 where $\hat{\mathbf{P}}_t = \hat{\mathbf{C}}_t / \sum_{a=0}^A \hat{C}_{a,t}$ is the expected age-composition proportion in year t .

245 Unless otherwise noted, the three EMs were correctly specified (fixed at the true values) for all
 246 model parameters except unfished recruitment (R_0), annual recruitment (R_t), parametric selectivity
 247 (S_{slope} and S_{50}), selectivity deviations ($\varepsilon_{a,t}$ in EMs *IID deviations* and *AR deviations*), and annual
 248 fully-selected fishing mortality (F_t). Among those estimated parameters, R_0 , S_{slope} , S_{50} and F_t were
 249 estimated as fixed effects, and $\varepsilon_{a,t}$ and R_t were estimated as random effects. We used Template
 250 Model Builder (TMB, Kristensen et al. (2016)) to implement mixed-effect parameter estimation.

251 In TMB, the marginal likelihood of fixed effect parameters was calculated using the Laplace
252 approximation to integrate across random effects (Kristensen et al. 2016) and fixed effect
253 parameters were then estimated via maximizing the marginal likelihood within the R (version 3.4.0)
254 computing environment (R Core Team 2017). We used the *nlm* function to minimize the
255 negative of the marginal log-likelihood and confirmed model convergence based on the
256 convergence flag the function provided and a positive-definite Hessian.

257 2.2. Simulation experiments

258 In this study, we conducted four related simulation experiments. A summary of the factorial
259 design of the OM and EM in each experiment can be found in Table 3.

260 2.2.1. What is the impact of under, right, or over-weighting on model performance?

261 The first simulation experiment aimed to evaluate the sensitivity of estimation performance of
262 the three EMs to data weighting, given that the true fishery selectivity had independent or
263 autocorrelated deviations. We compared the performance of each EM in estimating SB under three
264 data weighting scenarios: 1) under-weighting age-composition data by a factor of 10, which was
265 realized by setting $n_{input} = 0.1 \times n_{true} = 20$ in the three EMs; 2) right-weighting age-
266 composition data, which was realized by setting $n_{input} = n_{true} = 200$ in the three EMs; and 3)
267 over-weighting age-composition data by a factor of 10, which was realized by setting $n_{input} =$
268 $10 \times n_{true} = 2000$ in the three EMs. In this simulation experiment, four hundred simulation
269 replicates with unique process errors (in recruitment and selectivity) and observation errors (in
270 abundance index and age-composition observations) were generated for every combination of
271 population dynamics and OM case. Each EM (*zero deviations*, *IID deviations*, or *AR deviations*)
272 was then fitted to every generated simulation replicate individually under three data weighting

273 scenarios (under, right, or over-weighting). We evaluated the estimation performance of the three
 274 EMs based on the mean absolute relative error (MARE) in the estimate of final year SB:
 275 $\text{mean}(|\widehat{SB}_{t=20}/SB_{t=20} - 1|)$. This metric took both accuracy and precision into consideration.

276 *2.2.2. How well can we estimate effective sample size given time-varying selectivity*

277 The second simulation experiment aimed to compare three widely-used data weighting
 278 methods in stock assessments:

- 279 • McAllister-Ianelli (M-I) method (McAllister and Ianelli 1997). The effective sample size for
 280 the multinomial distribution was iteratively estimated through a tuning algorithm. In this study,
 281 it was computed as the harmonic mean of annual effective sample sizes

282
$$n_{eff} = \frac{T}{\sum_{t=1}^T \left(\frac{1}{n_{eff_t}} \right)} \quad (5a)$$

283
$$n_{eff_t} = \frac{\sum_{a=0}^A (\hat{P}_{a,t}(1 - \hat{P}_{a,t}))}{\sum_{a=0}^A (P_{a,t} - \hat{P}_{a,t})^2} \quad (5b)$$

284 and iteratively tuned until its relative difference between two iterations was less than 5%.

- 285 • Francis method (Francis 2011). The effective sample size for the multinomial distribution was
 286 also iteratively estimated through a tuning algorithm. Specifically, it is the inverse of the
 287 variance of normalized differences between the observed (P'_t) and expected (\hat{P}'_t) mean ages in
 288 age-composition

289
$$n_{eff} = \frac{1}{\text{Var} \left(\frac{P'_t - \hat{P}'_t}{\sqrt{v_t}} \right)} \quad (6a)$$

290
$$P'_t = \sum_{a=0}^A (aP_{a,t}) \quad (6b)$$

291
$$\hat{P}'_t = \sum_{a=0}^A (a\hat{P}_{a,t}) \quad (6c)$$

292
$$v_t = \sum_{a=0}^A (a^2\hat{P}_{a,t}) - \hat{P}'_t{}^2 \quad (6d)$$

293 and iteratively tuned until its relative difference between two iterations was less than 5%.

- 294 • Dirichlet-multinomial (D-M) method (Thorson et al. 2017). Different from the two tuning
 295 methods above, the D-M method estimated the effective sample size based on maximum
 296 likelihood. By assuming that age-composition data followed the linear parameterization of the
 297 Dirichlet-multinomial distribution, the effective sample size of the age-composition data was
 298 computed as

299
$$n_{eff} = \frac{1 + \theta n_{input}}{1 + \theta} \quad (7)$$

300 where n_{input} was the input sample size of the age-composition data. The D-M method
 301 estimated the effective sample size by fixing age-composition data ($n_{input}P_{a,t}$) and instead
 302 estimating θ as a parameter. The likelihood associated with the age-composition data was

303
$$L_{comp} \propto \prod_{t=1}^T \left(\frac{\Gamma(\theta n_{input})}{\Gamma(n_{input} + \theta n_{input})} \prod_{a=0}^A \frac{\Gamma(n_{input}P_{a,t} + \theta n_{input}\hat{P}_{a,t})}{\Gamma(\theta n_{input}\hat{P}_{a,t})} \right) \quad (8)$$

304 The three data weighting methods were compared based on two metrics: the ratio of estimated
 305 effective sample size to true sample size (n_{eff}/n_{true}) and the MARE in the estimate of final year
 306 SB. In this simulation experiment, four hundred simulation replicates with unique process errors
 307 and observation errors were generated for every combination of population dynamics and OM case.

308 OM case *Independent* approximately matched the simulation scenario explored in Thorson et al.
309 (2017), but other OM cases represented the first attempt to explore the sensitivity of the D-M
310 method to model mis-specification.

311 2.2.3. How does over-dispersion affect D-M estimates given time-varying selectivity?

312 In the third simulation experiment, we evaluated the performance of the three data weighting
313 methods in estimating the effective sample size of over-dispersed age-composition data. Over-
314 dispersed age-composition data were simulated by assuming that the extent of over-dispersion is
315 constant across age and time:

$$316 \quad \tilde{A}_{a,t} = A_{a,t} \times d \quad (9)$$

317 where $d (>1)$ denotes the extent of over-dispersion in age-composition data. For the two tuning
318 methods, the estimated effective sample size was a function of age-composition proportion ($P_{a,t}$),
319 which, under this assumption, did not change with the extent of over-dispersion in age-composition
320 data ($P_{a,t} = \tilde{A}_{a,t} / \sum_{a=0}^A \tilde{A}_{a,t} = A_{a,t} \times d / \sum_{a=0}^A (A_{a,t} \times d) = A_{a,t} / \sum_{a=0}^A A_{a,t}$). Thus, the estimated
321 effective sample size based on either tuning method should not be affected by the over-dispersion
322 in age-composition data. The focus of this simulation experiment was indeed on the D-M method,
323 for which n_{input} was specified to be the actual sample size (number of fish sampled; $n_{true} \times d$)
324 of the over-dispersed age-composition data. Eq. 9 simulated a type of over-dispersion case that all
325 fish were caught in groups of d individuals with identical age. By this definition, the true sample
326 size would be n_{true} .

327 We computed the ratio of estimated effective sample size to true sample size (n_{eff}/n_{true}) for
328 evaluating the performance of the D-M method with respect to estimating the effective sample
329 size, given that the age-composition data are over-dispersed. Due to the high computation demand

330 in this simulation experiment, we generated one hundred simulation replicates with unique process
331 errors and observation errors for every combination of population dynamics and OM case.

332 2.2.4. Can we estimate time-varying selectivity penalty and composition weighting simultaneously?

333 Lastly, we considered a more realistic situation where the degree of variation in selectivity was
334 estimated rather than known without error. In this simulation experiment, the degree of variation
335 in selectivity was iteratively estimated using a tuning algorithm inspired by Methot and Taylor
336 (2011) and introduced by Xu et al. (2018)

$$337 \quad \hat{\sigma}_S^2 = \text{SD}(\hat{\boldsymbol{\epsilon}})^2 + \frac{1}{6T} \sum_{a=2}^7 \sum_{t=1}^T \text{SE}(\hat{\epsilon}_{a,t})^2 \quad (10)$$

338 To replicate the case of Stock Synthesis (Methot and Wetzel 2013) and other widely-used
339 penalized likelihood models, here $\hat{\sigma}_S$ was estimated via the tuning approach instead of the mixed-
340 effect approach (i.e., EM3 instead of EM4 in Xu et al. 2018). Xu et al. (2018) showed that this
341 tuning algorithm could accurately estimate $\hat{\sigma}_S$ when the effective sample size was known without
342 error. In real-world assessments, however, both n_{eff} and $\hat{\sigma}_S$ are unknown and need to be estimated.

343 The focus of this simulation experiment was on the combined performance of the D-M method
344 for estimating n_{eff} and the tuning method for estimating $\hat{\sigma}_S$. How the effective sample size and
345 selectivity hyper-parameters were simultaneously estimated in this simulation are described below:

- 346 • Step 1: Tune selectivity variability ($\hat{\sigma}_S$) and effective sample size (n_{eff}). $\hat{\sigma}_S$ was iteratively
347 tuned in EM *IID deviations* until matching Eq. 10 within an accuracy of 0.01 while the D-M
348 method was used in every iteration of $\hat{\sigma}_S$ to estimate n_{eff} . $\hat{\sigma}_S$ was then fixed in EM *IID*
349 *deviation* and the estimated selectivity deviations were extracted.

350 • Step 2: Estimate selectivity autocorrelations ($\hat{\rho}_a$ and $\hat{\rho}_t$). $\hat{\rho}_a$ and $\hat{\rho}_t$ were estimated based on
351 an “external” estimation method (for more details see the description of EM *AR deviations* in
352 Xu et al. 2018). In brief, the two autocorrelation coefficients were estimated using the
353 maximum likelihood approach by fitting an external stand-alone model to selectivity
354 deviations that EM *AR deviations* estimated in step 1. The external stand-alone model
355 estimated $\hat{\rho}_a$ and $\hat{\rho}_t$ by assuming that selectivity deviations follow the multivariate normal
356 distribution described in Eq. 2 and that both $\hat{\rho}_a$ and $\hat{\rho}_t$ are between 0 and 1 (realized by using
357 the logit transformation).

358 The combined performance was evaluated according to the ratios of estimated to true values
359 of both n_{eff} and $\hat{\sigma}_S$. In addition, we compared the two estimated autocorrelation coefficients with
360 the corresponding true values to evaluate the performance of the “external” estimation method for
361 selectivity autocorrelations. Due to the high computation demand in this simulation experiment,
362 we generated one hundred simulation replicates with unique process errors and observation errors
363 for every combination of species and OM case.

364

365 **3. Results**

366 *3.1. What is the impact of under, right, or over-weighting on model performance?*

367 Overall, over-weighting age-composition data generally performed worse than under-
368 weighting age-composition data to the same extent. Results of the first simulation showed that
369 over-weighting tended to cause a larger estimation error in the final year SB (Fig. 3) in comparison
370 to under-weighting. Results also showed that over-weighting consistently corresponded to

371 significantly worse estimation performance than under-weighting for EM *AR deviations*, the EM
372 with correctly-specified selectivity (Fig. 3; see Fig. A1 for resampled uncertainty levels).

373 Whether under-weighting or over-weighting age-composition data performed better varied
374 somewhat among species and OMs. Under OM *Independent*, right-weighting and over-weighting
375 age-composition data performed best and worse, respectively, regardless of how selectivity was
376 specified in the EM (Fig. 3, first column). Over-weighting also performed worse under OM *Time-*
377 *correlated*, except for EM *IID deviations* for Pacific hake, which performed worst when age-
378 composition data were down-weighted (Fig. 3, second column). It is worth noting that under-
379 weighting could produce the best performing EM (i.e., EM *IID deviations*) in this case. Under OM
380 *Age-correlated*, over-weighting performed better and worse than under-weighting when the
381 variation in selectivity was ignored (EM *zero deviations*) and estimated (EMs *IID deviations* and
382 *AR deviations*), respectively (Fig. 3, third column). Again, right-weighting generally performed
383 best regardless of how selectivity was specified in the EM. Under OM *Age- and time-correlated*,
384 the three data weighting scenarios performed similarly for EMs *zero deviations* and *IID deviations*,
385 at least within the weighting range ($0.1\times$ - $10\times$) investigated in this study (Fig. 3, fourth column).
386 For EM *AR deviations*, over-weighting and right-weighting consistently performed worst and best,
387 respectively.

388 3.2. How well can we estimate the effective sample size given time-varying selectivity?

389 Among the three data weighting methods (M-I, Francis, and D-M methods), the D-M method
390 provided the most accurate estimated effective sample size regardless of whether the EM allowed
391 for time-varying selectivity. For EM *zero deviations* which mistakenly specified constant
392 selectivity, all three data weighting methods estimated a reduced effective sample size (medians
393 within $0.2\times$ - $0.7\times$ of the true sample size) (Fig. 4). This behavior was expected given that these

394 models are mis-specified. For EMs *IID deviations* and *AR deviations*, the effective sample size
395 that the D-M method provided was very accurate while those the M-I method and especially the
396 Francis method estimated were considerably larger than the true sample size (Fig. 4). From a
397 median point of view, the M-I method over-estimated the effective sample size by as much as 4×
398 and 13× for Pacific hake and Pacific sardine, respectively; the Francis method over-estimated the
399 effective sample size by as much as 8× and 20× for Pacific hake and Pacific sardine, respectively.
400 It is worth noting that, by definition (Eq. 7), the effective sample size estimated by the D-M method
401 must be smaller than or approximately the same as (when θ is estimated to be very large) the input
402 sample size, which in this simulation experiment is identical to the true sample size. Therefore, the
403 effective sample sizes estimated by the D-M method were all smaller than the true sample size in
404 this simulation experiment, regardless of whether the EM accounts for time-varying selectivity or
405 not.

406 Because the D-M method provided the most accurate estimated effective sample size when
407 allowing for time-varying selectivity, both EMs *IID deviations* and *AR deviations* performed best
408 when using the D-M method for data weighting. The first simulation informed us that estimation
409 performance was relatively insensitive to data weighting when age-composition data were under-
410 weighted (Fig. 4). Since all three data weighting methods under-estimated the effective sample
411 size under constant selectivity, the method for data weighting had little effect on the estimation
412 performance of EM *zero deviations* (Fig. 5). For EMs *IID deviations* and *AR deviations*, in
413 comparison, estimation performance could be significantly affected by the method on which data
414 weighting was based: the D-M and Francis method generally corresponded to smallest and largest
415 errors in the estimate of final year SB, respectively (Fig. 5). That was because, when allowing for
416 time-varying selectivity, the M-I method and especially the Francis method considerably over-

417 estimated the effective sample size, which was, in contrast, slightly under-estimated by the D-M
418 method (Fig. 4). Also, the extent to which the two tuning methods over-estimated the effective
419 sample size when estimating time-varying selectivity was larger for Pacific sardine than Pacific
420 hake. Consequently, the benefit of the D-M method in terms of improving the estimate of final
421 year SB was significant for Pacific sardine, but not for Pacific hake (Fig. 5; see Fig. A2 for
422 resampled uncertainty levels).

423 Some cases in this simulation suggested that estimating time-varying selectivity (EMs *IID*
424 *deviations* and *AR deviations*) resulted in worse performance than assuming time-invariant
425 selectivity (EM *zero deviations*) when using either the Francis or M-I method for data weighting
426 (Fig. 5). This is surprising, given that the true selectivity in the OM was simulated to have a
427 moderate level of variation over both age and time (see Eq. 1 and Fig. 1). This pattern indicated
428 that improving stock assessments by accounting for time-varying selectivity requires
429 simultaneously addressing data weighting and the time-varying process.

430 3.3. How does over-dispersion affect D-M estimates given time-varying selectivity?

431 The effective sample size of over-dispersed age-composition data was under-estimated by the
432 D-M method for all three EMs, but the degree of under-estimation in all cases remained in a
433 reasonable range (medians larger than $0.2\times$ of the true sample size) (Fig. 6). As expected, the
434 effective sample size was estimated to be considerably below the input sample size (medians
435 within $0.2\times$ - $0.6\times$ of the true sample size) for EM *zero deviations*. This result is expected given that
436 this EM is mis-specified. When allowing for time-varying selectivity (EMs *IID deviations* and *AR*
437 *deviations*), the median estimated effective sample size that the D-M method provided was above
438 $0.5\times$ of the true sample size under all degrees of over-dispersion (Fig. 6). The bias in the estimated
439 effective sample size became slightly greater as the degree of over-dispersion increased from 2 to

440 10, which was in accordance with the trend found in a previous study (see Fig. 4 in Thorson et al.
441 (2017). We also noted that although the difference was not dramatic, the D-M method generally
442 performed better for Pacific hake simulations than Pacific sardine simulations.

443 3.4. *Can we estimate time-varying selectivity penalty and composition weighting simultaneously?*

444 Like in the previous simulation experiment where the level of variation in selectivity was
445 assumed known without error, the median estimated effective sample size that the D-M method
446 provided was still above $0.5\times$ of the true sample size (Fig. 7). Moreover, the bias in the estimated
447 effective sample size was still positively related to the degree of over-dispersion in age-
448 composition data. The first simulation suggested that the performance of models that estimated
449 time-varying selectivity was not sensitive to down-weighting age-composition data. As such, it
450 was not surprising to find that MARE was negligibly impacted by the under-estimation of the
451 effective sample size (Fig. A3) within this range of degree of under-estimation (by 10%-50%) (Fig.
452 7). Namely, model performance was not sensitive to the degree of over-dispersion in age-
453 composition data when using the D-M method for data weighting. Importantly, in combination
454 with the D-M method for data weighting, the tuning method that was developed by Xu et al. (2018)
455 was useful for estimating the level of variation in selectivity. The level of variation in selectivity
456 was only slightly under-estimated (medians within $0.75\times$ - $1\times$ of the true level), within which the
457 largest degree of under-estimation occurred under OM *Age- and time-correlated*.

458 The “external” estimation method for the two autocorrelation coefficients for selectivity
459 deviations ($\hat{\rho}_a$ and $\hat{\rho}_t$) were also useful (Fig. 8). Under OM *Independent* where the true
460 coefficients were both 0, the median estimates that the “external” method provided were between
461 0 and 0.2. Under OMs *Age-dependent* and *Time-dependent*, where one of the two true coefficients
462 was positive (0.8), the median estimate of that coefficient was only slightly under-estimated

463 (median larger than 0.6). Under OM *Age- and time-dependent* where the true coefficients were
464 both positive (0.8), the median estimates of the two coefficients were mostly above 0.4 In
465 accordance with the finding in Xu et al. (2018), this “external” estimation method generally
466 performed better for $\hat{\rho}_t$ than $\hat{\rho}_a$ (Fig. 8).

467

468 **4. Discussion**

469 This study aimed to compare the performance of widely-used data weighting methods in the
470 assessment models that allow for time-varying selectivity. We conducted four simulation
471 experiments to evaluate the sensitivity of model estimates to data weighting, and more importantly,
472 to evaluate the performance of M-I, Francis, and D-M methods given time-varying selectivity. For
473 assessment models that estimated time-varying selectivity, over-weighting generally led to larger
474 estimation error in final year SB than did under-weighting to the same extent, suggesting that over-
475 weighting should be avoided even when allowing for time-varying selectivity. Among the three
476 data weighting methods compared in this study, the D-M method out-performed the other two
477 methods when estimating time-varying selectivity. Moreover, the D-M method was still useful
478 even when age-composition data were over-dispersed and the level of variation in selectivity was
479 simultaneously estimated. In conclusion, the D-M method was recommended over the M-I and
480 Francis methods for the assessments that explore time-varying selectivity.

481 Overall, over-weighting composition data tends to cause larger estimation error in final year
482 SB than does under-weighting composition data to the same extent. The result echoes Francis’
483 (2011) suggestion that age-composition data should not be over-weighted. This suggestion was
484 made based on the idea that while composition data are important to inform selectivity and

485 recruitment variation, the estimated population trend should be primarily driven by the more
486 reliable abundance indices, especially when data conflict exists between abundance indices and
487 composition data. Importantly, this study shows that estimation performance is more degraded by
488 over-weighting than under-weighting regardless of whether selectivity is mis-specified or not. In
489 some cases, data weighting had a larger impact on the estimation performance of assessment
490 models with correctly-specified selectivity than those without, implying that data weighting is
491 important for stock assessments with any selectivity specifications.

492 For assessment models that estimate time-varying selectivity, the D-M method overall
493 performs better than the M-I and Francis methods with regards to weighting age-composition data.
494 Under the specification of time-varying selectivity, the M-I method and especially the Francis
495 method over-estimate the effective sample size greatly and consequently correspond to large
496 estimation error in final year SB. The fact that the effective sample size is greatly over-estimated
497 by the two tuning methods is likely due to the expected and observed age-compositions tend to
498 match more closely under a more flexible (i.e., time-varying) selectivity specification (Francis
499 2017, Punt et al. 2014). In contrast, the D-M method under-estimates the effective sample size
500 slightly and consequently corresponds to smaller estimation error in final year SB. When age-
501 composition data are over-dispersed, simulation results show that the D-M method under-estimates
502 the effective sample size to a certain extent. However, the extent of the under-estimation is still
503 smaller in comparison with the extent to which the two tuning methods over-estimate the effective
504 sample size of randomly-sampled age-composition data. Furthermore, assessment models that use
505 the D-M data weighting method consider the uncertainty about data weighting. the D-M method
506 estimates effective sample size as a parameter of the assessment model based on maximum
507 likelihood, so it can propagate the uncertainty about data weighting into the confidence interval of

508 estimated population and management attributes (Thorson et al. 2017). In contrast, both the two
509 tuning methods ignore the uncertainty in estimated effective sample size, leading to under-
510 estimated uncertainty in model estimates (Maunder 2011). However, it should be noted that when
511 using the D-M method for data weighting, the confidence interval of estimated population and
512 management attributes could be over-estimated as the simulations in this study show that the D-M
513 method tends to under-estimate the effective sample size.

514 For assessment models that estimate time-varying selectivity, the D-M data weighting method
515 is robust for over-dispersed age-composition data, which is a common phenomenon in fisheries.
516 Within the range of over-dispersion we investigated ($2\times$ - $10\times$), the D-M method under-estimated
517 the effective sample size by less than 50%, regardless of the OM case and type of life history. The
518 first simulation experiment showed that under-weighting age-composition data by such an extent
519 should only minorly degrade the estimation performance of the assessment model. It should be
520 noted that the D-M method tends to down-weight age-composition data to a much larger extent
521 when selectivity is specified to be constant than time-varying, implying that data weighting based
522 on the D-M method is informed by the goodness-of-fit of age-composition data in the assessment
523 model (Thorson et al. 2017). In addition to the linear parameterization used in the paper, the
524 Dirichlet-multinomial can also be parameterized in another way (i.e., the saturation
525 parameterization). In the simulations in the paper, the input sample size is specified to be identical
526 among years, in which case the two parameterizations result in identical parameter estimates.
527 Future research could compare the two parameterizations when input sample size varies among
528 years.

529 We are aware that the comparison of the three data weighting methods in the second simulation
530 experiment is tilted towards the D-M method. By construction, the D-M method only allows the

531 effective sample size to be smaller than the input sample size (which is specified to be the true
532 sample size in that simulation), leading to a more restricted parameter space for the effective
533 sample size and a larger probability for the effective sample size to be close to the true sample size.
534 To make a fairer comparison between the three data-weighting methods, we then conducted the
535 third simulation experiment in which the effective sample size can be as large as 10x of the true
536 sample size. Being consistent with the pattern found in Thorson et al. (2017), the performance of
537 the D-M method is negatively related to the ratio of input sample size to true sample size. However,
538 even when the input sample size for the D-M method is specified to be 10x of the true sample size,
539 the D-M method still performs better than the two tuning methods. We recommend ongoing
540 research to accurately estimate the input sample size from field-measurements of age- and length-
541 composition (Stewart and Hamel 2014, Thorson 2014, Thorson and Haltuch 2018) because an
542 accurate starting point for weighting compositional data improves model performance when using
543 the D-M data weighting method.

544 In terms of estimation performance of an assessment model, correctly specifying the
545 distributional penalty for selectivity deviations is as important as choosing a proper method (i.e.,
546 the D-M method) for data weighting. When using the D-M data weighting method, correctly
547 specifies selectivity (EM *AR deviations*) greatly outperformed the other two EMs with mis-
548 specified selectivity (EMs *zero deviations* and *IID deviations*). Several studies have suggested
549 considering data weighting and time-varying selectivity together in stock assessments (Francis
550 2011, Stewart and Monnahan 2017, Thorson et al. 2017, Wang and Maunder 2017). Results from
551 this study provide another strong support for this suggestion.

552 According to our simulation study, the D-M method for weighting composition data and the
553 tuning method for penalizing selectivity variation (Xu et al. 2018) are able to provide proper data

554 weighting and selectivity penalizing simultaneously. In real-world stock assessments, both the
555 level of variation in selectivity and the level of over-dispersion in composition data are unknown
556 and need to be estimated. Considering that both methods have been implemented in Stock
557 Synthesis (Methot and Wetzel 2013), a widely used stock assessment package, we recommend
558 users to explore the two methods together in real-world stock assessments. When exploring the
559 two methods simultaneously in stock assessments, we also recommend evaluating the
560 autocorrelations in selectivity deviations using the “external” estimation method, which performs
561 reasonably well in Xu et al. (2018) as well as in this study. It can improve stock assessments if
562 process errors in selectivity are highly autocorrelated.

563 We note that the performance of the D-M method is likely over-estimated in our idealized
564 simulations. First, our assumption about selectivity deviations in the OM allows the simulated
565 time-varying selectivity to be asymptotic only. However, real fishery selectivity can be dome-
566 shaped (Sampson and Scott 2011, Waterhouse et al. 2014) and the multinomial distribution was
567 found to perform poorly in simulations where selectivity is assumed to be dome-shaped. Second,
568 we only evaluated the impacts of mis-specifying selectivity on the performance of the D-M method
569 in this study. Other biological processes including natural mortality, growth, and maturity were all
570 assumed known without error. In real-world stock assessments, however, these biological
571 processes are likely to vary in complicated ways, such that assessment models are likely
572 misspecified in multiple ways simultaneously. In other words, these biological processes are more
573 or less mis-specified in real-world stock assessments, leading to larger discrepancies between
574 observed and predicted age-composition. Considering that the D-M data weighting method already
575 under-estimates the effective sample size in this study, it may under-estimate the effective sample
576 size to a larger extent in real-world stock assessments.

577 Third, there is another obvious limitation of this simulation study that can cause over-
578 estimating the performance of the D-M method. In this simulation study, the D-M method was
579 used to weight the age-composition data sampled using a closely-related distribution (i.e.,
580 multinomial). Studies (Berg and Nielsen 2016, Berg et al. 2014) have shown that sampling errors
581 in real fishery age-compositions can be positively correlated among ages. The multinomial
582 distribution, however, only allows negative correlations among ages and therefore may not be
583 appropriate for the sampling model that generates age-composition samples (Albertsen et al. 2017,
584 Francis 2014). Indeed, the D-M method may perform worse for length-composition data because
585 the positive correlations among lengths tend to be higher than those among ages. As such, the
586 performance of the D-M method needs to be more closely evaluated in future studies using real
587 fishery age-composition data.

588

589 **Acknowledgments**

590 We thank Mark Maunder, Jim Ianelli, the associated editor, and two anonymous reviewers for
591 providing valuable inputs to an earlier draft. We also thank Kasper Kristensen for developing
592 Template Model Builder. Haikun Xu was partially funded by NOAA Award NA15OAR4320063
593 through the University of Washington Joint Institute for the Study of the Atmosphere and Ocean
594 (JISAO).

595 **References**

- 596 Albertsen, C.M., Nielsen, A., and Thygesen, U.H. 2017. Choosing the observational likelihood in
597 state-space stock assessment models. *Canadian Journal of Fisheries and Aquatic Sciences*
598 **74(5): 779-789.**
- 599 Berg, C.W., and Nielsen, A. 2016. Accounting for correlated observations in an age-based state-
600 space stock assessment model. *ICES Journal of Marine Science: Journal du Conseil:*
601 *fsw046.*
- 602 Berg, C.W., Nielsen, A., and Kristensen, K. 2014. Evaluation of alternative age-based methods for
603 estimating relative abundance from survey data in relation to assessment models. *Fisheries*
604 *research* **151: 91-99.**
- 605 Francis, R.I.C.C. 2011. Data weighting in statistical fisheries stock assessment models. *Canadian*
606 *Journal of Fisheries and Aquatic Sciences* **68(6): 1124-1138.**
- 607 Francis, R.I.C.C. 2014. Replacing the multinomial in stock assessment models: A first step.
608 *Fisheries Research* **151: 70-84.**
- 609 Francis, R.I.C.C. 2017. Revisiting data weighting in fisheries stock assessment models. *Fisheries*
610 *Research* **192: 5-15.**
- 611 Hulson, P.-J.F., Hanselman, D.H., and Quinn, T.J. 2011. Determining effective sample size in
612 integrated age-structured assessment models. *ICES Journal of Marine Science* **69(2): 281-**
613 **292.**
- 614 Kotwicki, S., Horne, J.K., Punt, A.E., and Ianelli, J.N. 2015. Factors affecting the availability of
615 walleye pollock to acoustic and bottom trawl survey gear. *ICES Journal of Marine Science*
616 **72(5): 1425-1439.**

617 Kristensen, K., Nielsen, A., Berg, C., Skaug, H., and Bell, B. 2016. Template model builder TMB.
618 J. Stat. Softw **70**: 1-21.

619 Lowe, S.A., Ianelli, J.N., and Palsson, W. 2017. Assessment of the Atka mackerel stock in the
620 Bering Sea and Aleutian Islands [online]. Available from
621 www.afsc.noaa.gov/REFM/Docs/2012/BSAIatka.pdf.

622 Martell, S., and Stewart, I. 2014. Towards defining good practices for modeling time-varying
623 selectivity. Fisheries Research **158**: 84-95.

624 Maunder, M.N. 2011. Review and evaluation of likelihood functions for composition data in stock-
625 assessment models: Estimating the effective sample size. Fisheries Research **109**(2): 311-
626 319.

627 Maunder, M.N., Crone, P.R., Punt, A.E., Valero, J.L., and Semmens, B.X.J.F.R. 2017. Data
628 conflict and weighting, likelihood functions and process error. (192): 1-4.

629 Maunder, M.N., and Piner, K.R. 2017. Dealing with data conflicts in statistical inference of
630 population assessment models that integrate information from multiple diverse data sets.
631 Fisheries Research **192**: 16-27.

632 Maunder, M.N., and Punt, A.E. 2013. A review of integrated analysis in fisheries stock assessment.
633 Fisheries Research **142**: 61-74.

634 McAllister, M.K., and Ianelli, J.N. 1997. Bayesian stock assessment using catch-age data and the
635 sampling-importance resampling algorithm. Canadian Journal of Fisheries and Aquatic
636 Sciences **54**(2): 284-300.

637 Methot, R.D. 2009. Stock assessment: operational models in support of fisheries management. *In*
638 The future of fisheries science in North America. Springer. pp. 137-165.

639 Methot, R.D., and Taylor, I.G. 2011. Adjusting for bias due to variability of estimated recruitments
640 in fishery assessment models. *Canadian Journal of Fisheries and Aquatic Sciences* **68**(10):
641 1744-1760.

642 Methot, R.D., and Wetzel, C.R. 2013. Stock synthesis: a biological and statistical framework for
643 fish stock assessment and fishery management. *Fisheries Research* **142**: 86-99.

644 Minte-Vera, C.V., Maunder, M.N., Aires-da-Silva, A.M., Satoh, K., and Uosaki, K. 2017. Get the
645 biology right, or use size-composition data at your own risk. *Fisheries Research* **192**: 114-
646 125.

647 Punt, A.E., Hurtado-Ferro, F., and Whitten, A.R. 2014. Model selection for selectivity in fisheries
648 stock assessments. *Fisheries Research* **158**: 124-134.

649 R Core Team. 2017. R: A language and environment for statistical computing. R Foundation for
650 Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

651 Sampson, D.B., and Scott, R.D. 2011. A spatial model for fishery age-selection at the population
652 level. *Canadian Journal of Fisheries and Aquatic Sciences* **68**(6): 1077-1086.

653 Sampson, D.B., and Scott, R.D. 2012. An exploration of the shapes and stability of population-
654 selection curves. *Fish and Fisheries* **13**(1): 89-104.

655 Stewart, I.J., and Hamel, O.S. 2014. Bootstrapping of sample sizes for length-or age-composition
656 data used in stock assessments. *Canadian journal of fisheries and aquatic sciences* **71**(4):
657 581-588.

658 Stewart, I.J., and Martell, S.J. 2014. A historical review of selectivity approaches and retrospective
659 patterns in the Pacific halibut stock assessment. *Fisheries Research* **158**: 40-49.

660 Stewart, I.J., and Monnahan, C.C. 2017. Implications of process error in selectivity for approaches
661 to weighting compositional data in fisheries stock assessments. *Fisheries Research* **192**:
662 126-134.

663 Thorson, J.T. 2014. Standardizing compositional data for stock assessment. *ICES Journal of*
664 *Marine Science* **71**(5): 1117-1128.

665 Thorson, J.T. 2018. Perspective: Let's simplify stock assessment by replacing tuning algorithms
666 with statistics.

667 Thorson, J.T., Clarke, M.E., Steward, I.J., and Punt, A. 2013a. The implications of spatially
668 varying catchability on bottom trawl surveys of fish abundance: a proposed solution
669 involving underwater vehicles. *70*(2): 294-306.

670 Thorson, J.T., and Cope, J.M. 2015. Catch curve stock-reduction analysis: An alternative solution
671 to the catch equations. *Fisheries Research* **171**: 33-41.

672 Thorson, J.T., and Haltuch, M.A. 2018. Spatiotemporal analysis of compositional data: increased
673 precision and improved workflow using model-based inputs to stock assessment. *Canadian*
674 *Journal of Fisheries and Aquatic Sciences*(999): 1-14.

675 Thorson, J.T., Johnson, K.F., Methot, R.D., and Taylor, I.G. 2017. Model-based estimates of
676 effective sample size in stock assessment models using the Dirichlet-multinomial
677 distribution. *Fisheries Research* **192**: 84-93.

678 Thorson, J.T., Minto, C., Minto-Vera, C.V., Kleisner, K.M., and Longo, C. 2013b. A new role for
679 effort dynamics in the theory of harvested populations and data-poor stock assessment.
680 *Canadian Journal of Fisheries and Aquatic Sciences* **70**(12): 1829-1844.

681 Thorson, J.T., and Taylor, I.G. 2014. A comparison of parametric, semi-parametric, and non-
682 parametric approaches to selectivity in age-structured assessment models. *Fisheries*
683 *Research* **158**: 74-83.

684 Venables, W., and Ripley, B. 2002. *Modern applied statistics* (Fourth S., editor) New York.
685 Springer.

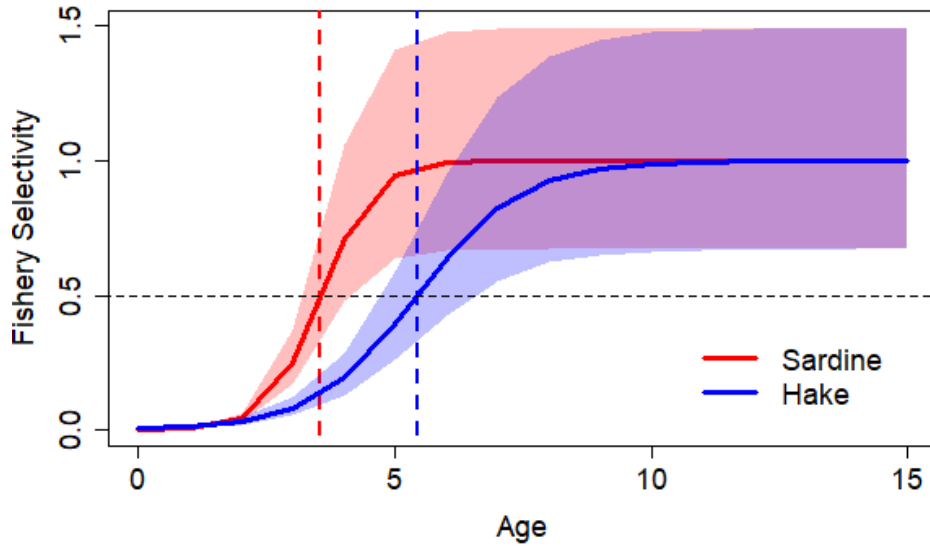
686 Wang, S.-P., and Maunder, M.N. 2017. Is down-weighting composition data adequate for dealing
687 with model misspecification, or do we need to fix the model? *Fisheries Research* **192**: 41-
688 51.

689 Waterhouse, L., Sampson, D.B., Maunder, M., and Semmens, B.X. 2014. Using areas-as-fleets
690 selectivity to model spatial fishing: asymptotic curves are unlikely under equilibrium
691 conditions. *Fisheries research* **158**: 15-25.

692 Xu, H., Thorson, J.T., Methot, R.D., and Taylor, I.G. 2018. A new semi-parametric method for
693 autocorrelated age-and time-varying selectivity in age-structured assessment models.
694 *Canadian Journal of Fisheries and Aquatic Sciences*(999): 1-18.

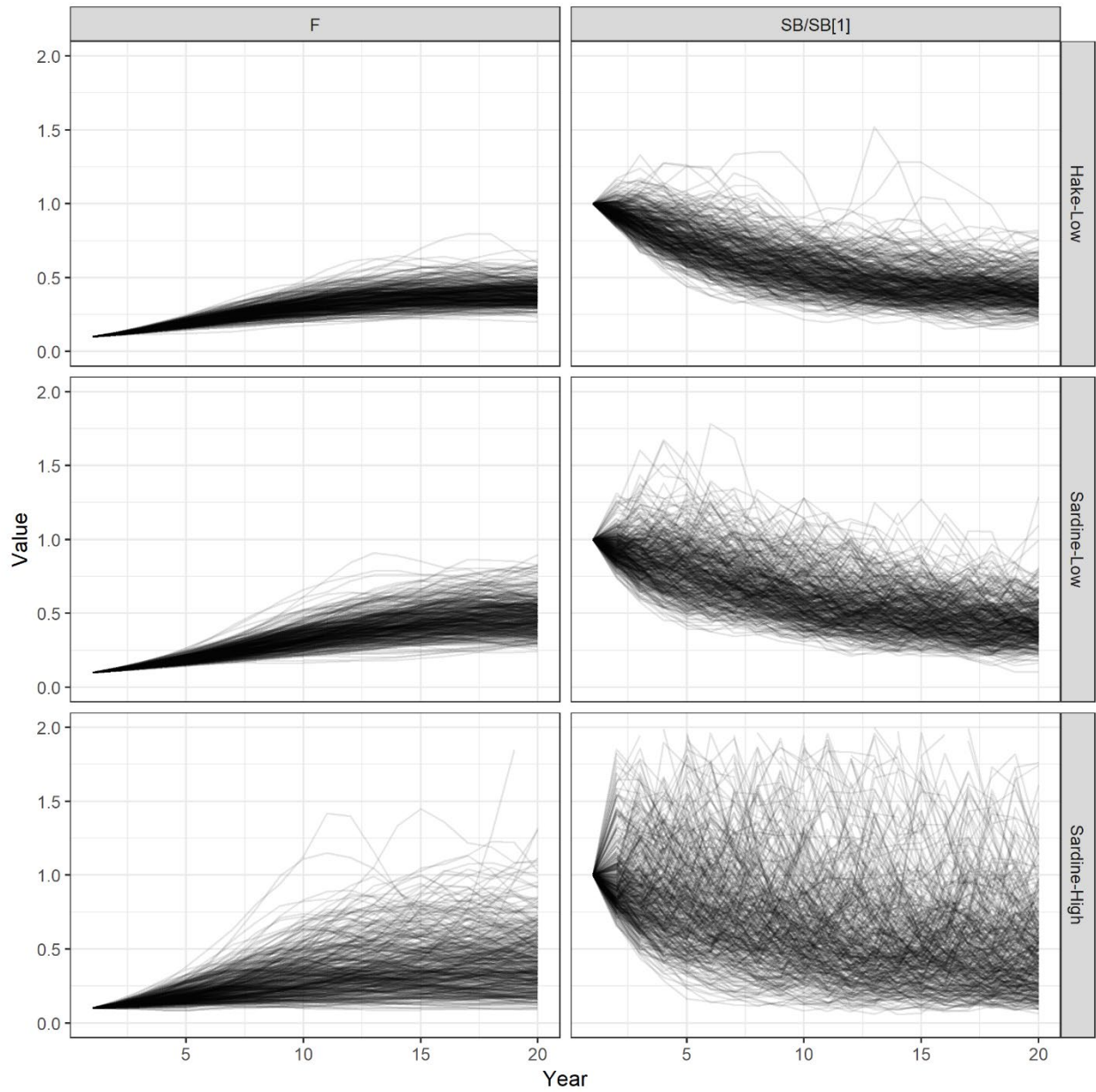
695

696 **Figures and Tables**

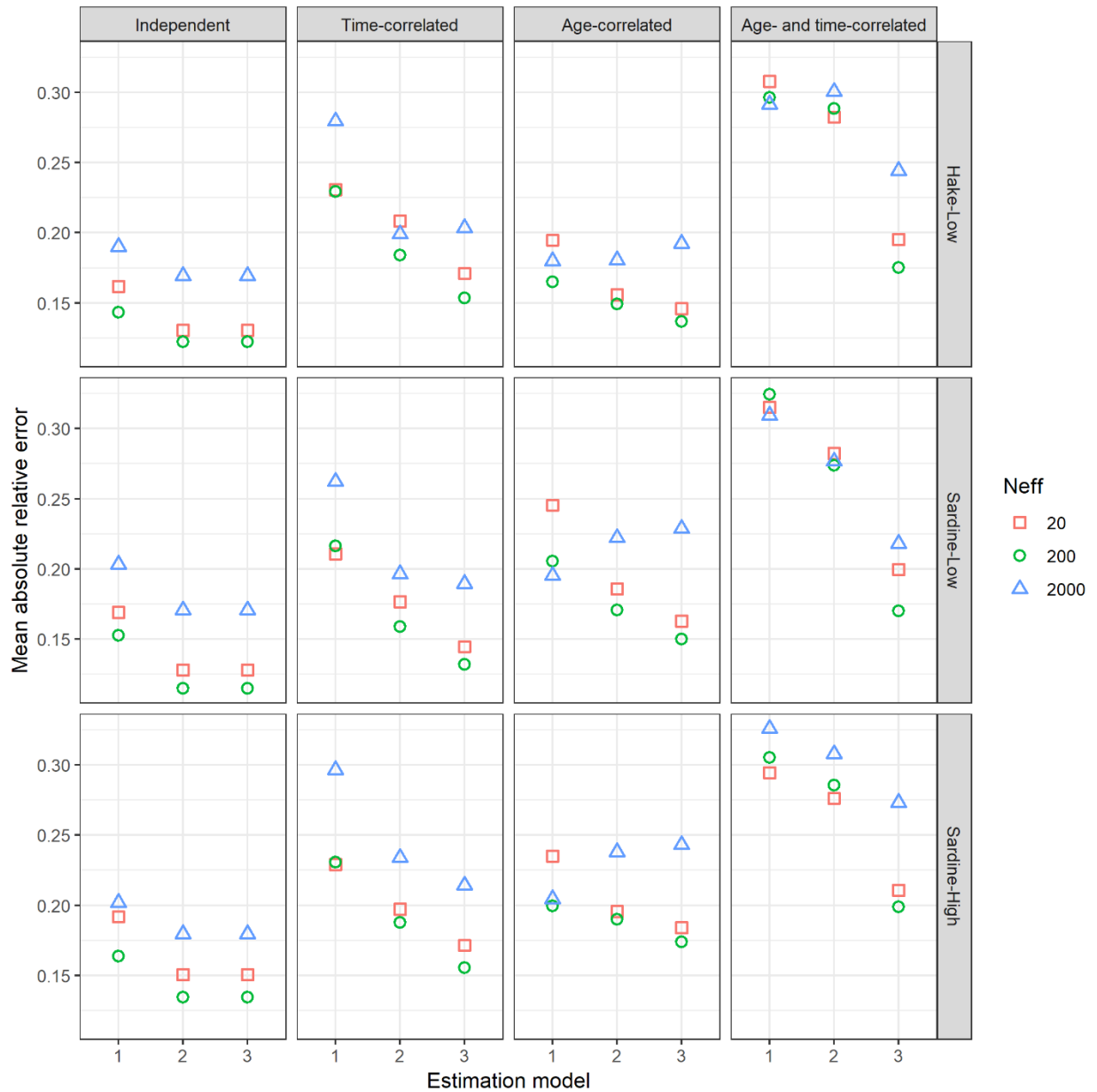


697

698 Figure 1. Comparison of the parametric fishery selectivity for the two types of life history (Pacific
699 hake and Pacific sardine) as a function of age. The shaded areas show the ± 1 standard deviation
700 range of selectivity variation. The vertical dashed lines mark the age at 50% selection of the fishery.



701
 702 Figure 2. 1st simulation experiment: trajectories of fully-selected fishing mortality (left) and
 703 spawning biomass (right) for the four hundred replicates. To facilitate the comparison among
 704 replicates, spawning biomass is rescaled to have an initial ($t = 1$) value of 1.



705

706 Figure 3. 1st simulation experiment: mean absolute relative error in the estimate of final year

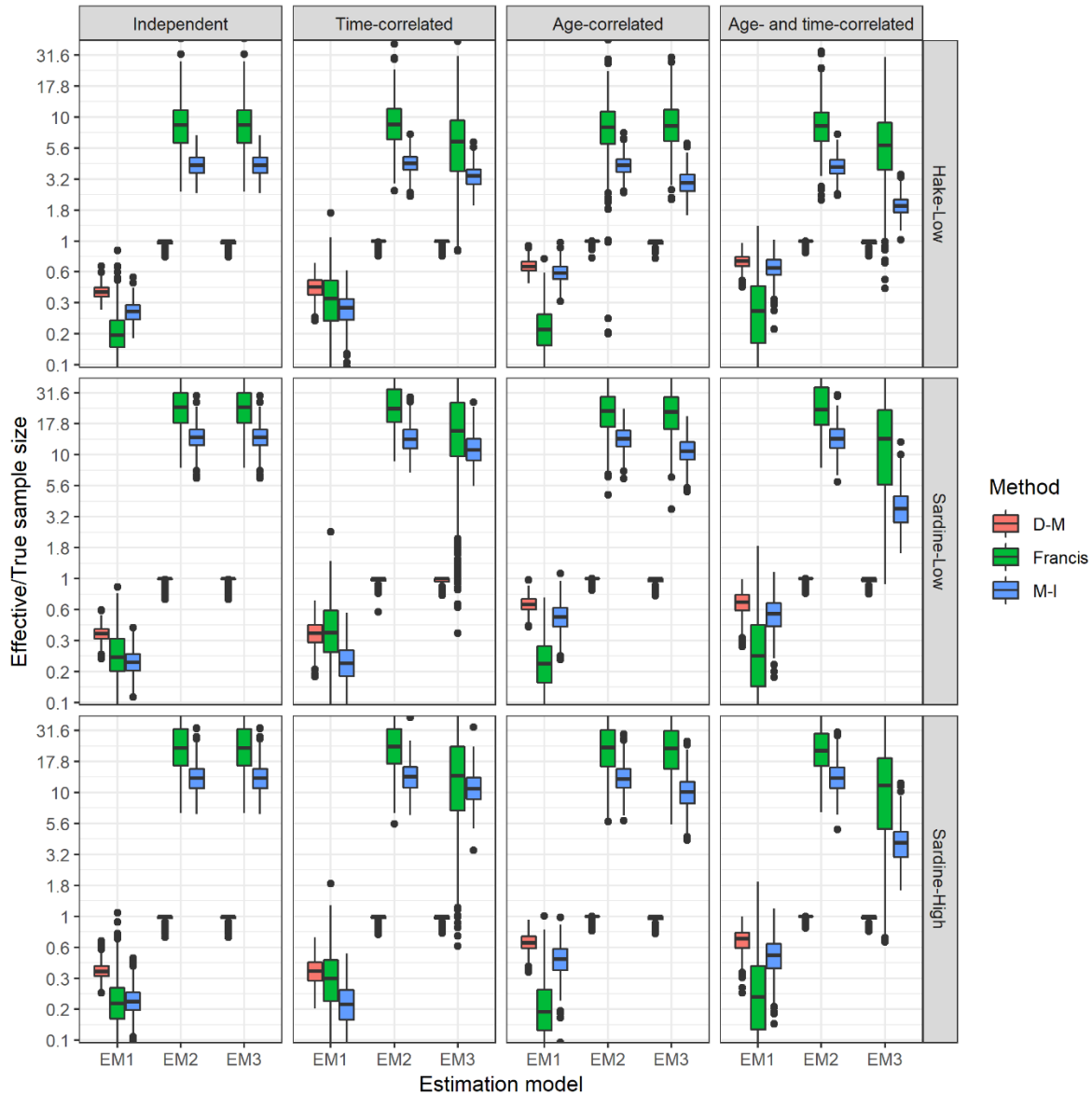
707 spawning biomass under the scenario of under-weighting (red circle), right-weighting (green

708 triangle), or over-weighting (blue square) age-composition data. The four columns correspond to

709 the four autocorrelation cases for simulated selectivity deviations: *Independent*, *Time-correlated*,

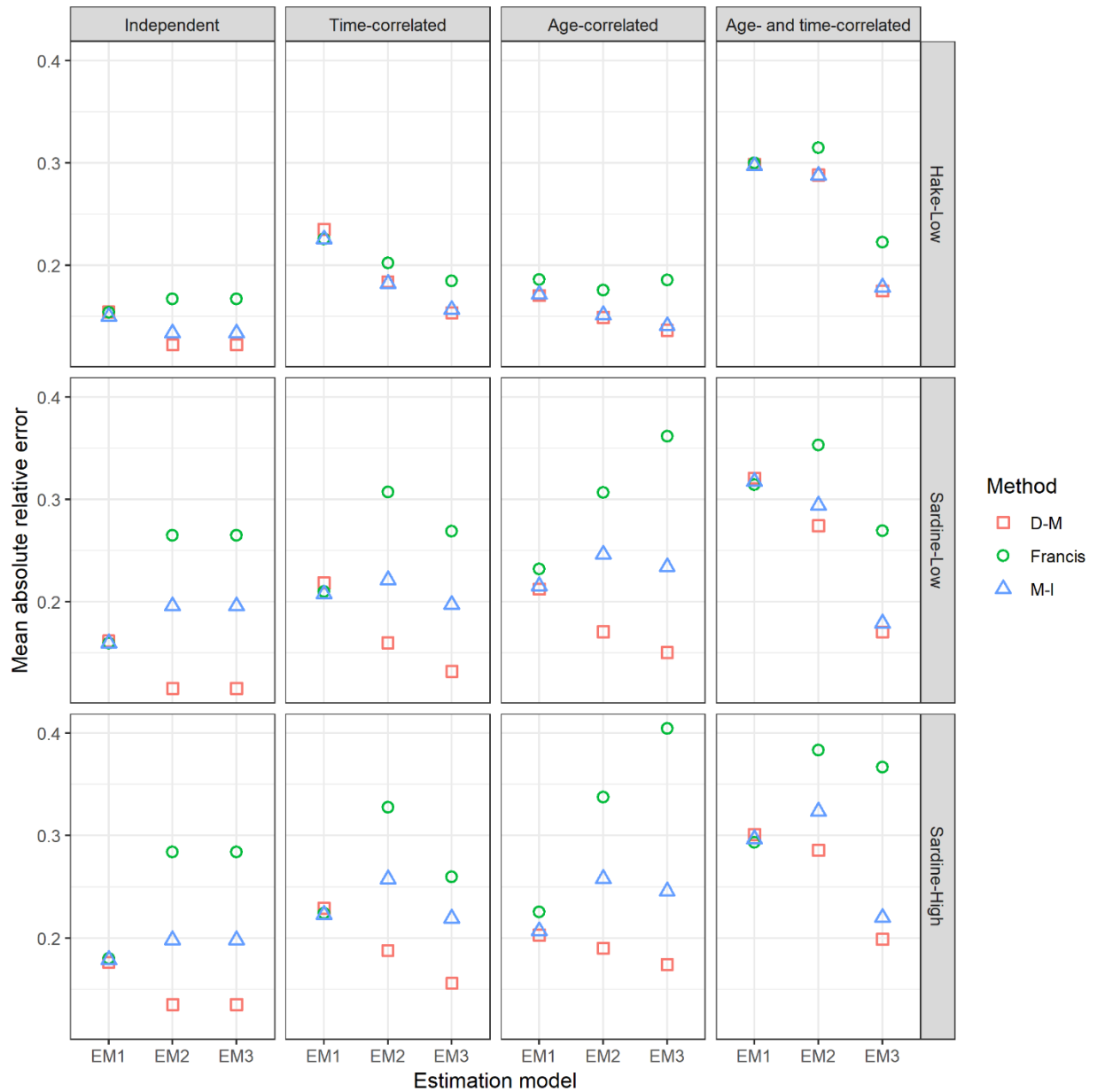
710 *Age-correlated*, and *Age- and time-correlated*. EM1-3 have different selectivity specifications:

711 *zero deviations*, *IID deviations*, and *AR deviations*.



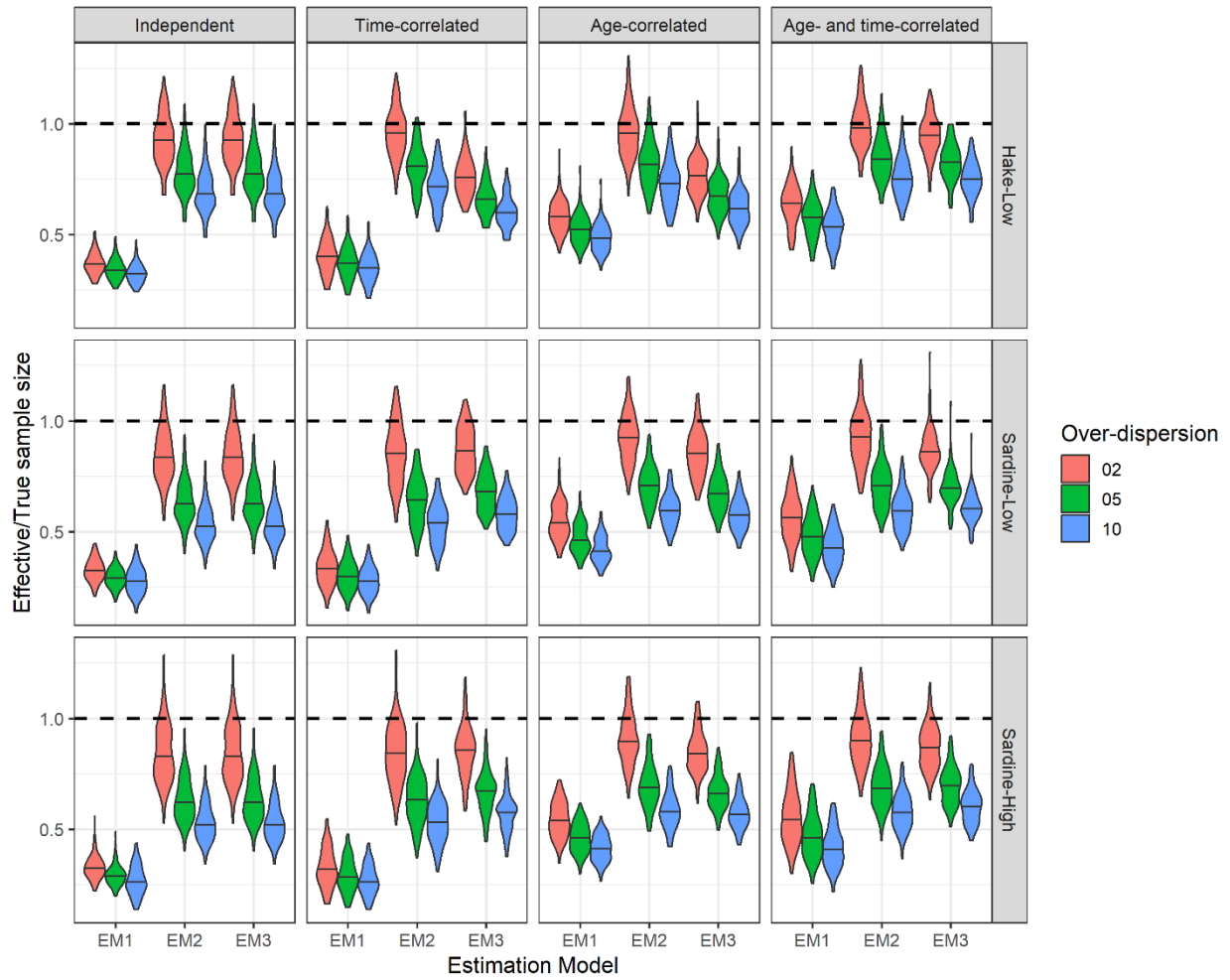
712

713 Figure 4. 2nd simulation experiment: boxplot for the ratio of effective sample size to true sample
 714 size using the Dirichlet-multinomial (D-M), Francis, and McAllister-Ianelli (M-I) methods. The
 715 lower and upper hinges mark the first and third quantiles and the two whiskers extend to the value
 716 no further than 1.5 interquartile range from the corresponding hinge. The four columns correspond
 717 to the four autocorrelation cases for simulated selectivity deviations: *Independent*, *Time-correlated*,
 718 *Age-correlated*, and *Age- and time-correlated*. EM1-3 have different selectivity specifications:
 719 *zero deviations*, *IID deviations*, and *AR deviations*.



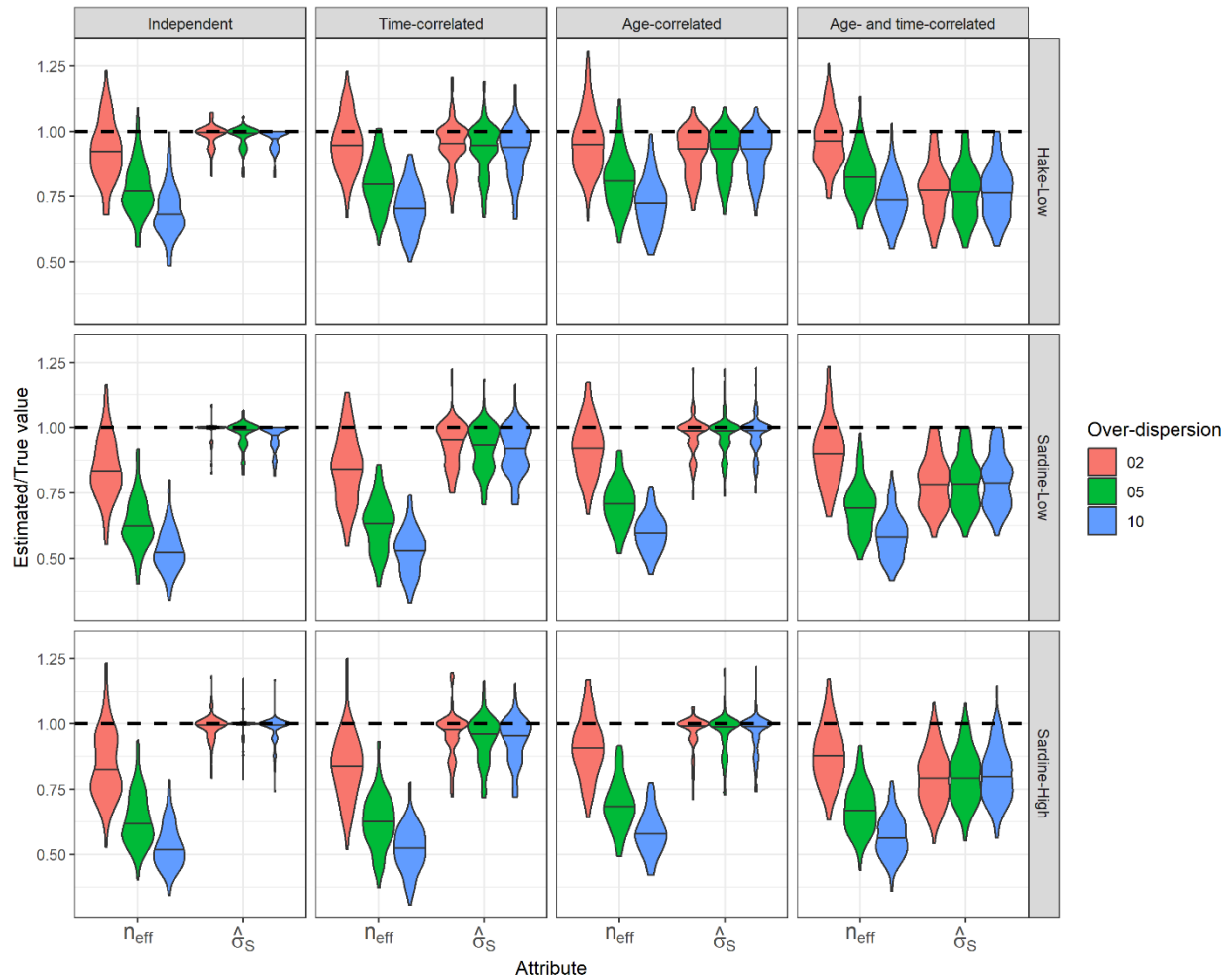
720

721 Figure 5. 2nd simulation experiment: mean absolute relative error in the estimate of final year
 722 spawning biomass under the Dirichlet-multinomial (D-M), Francis, and McAllister-Ianelli (M-I)
 723 methods. The four columns correspond to the four autocorrelation cases for simulated selectivity
 724 deviations: *Independent*, *Time-correlated*, *Age-correlated*, and *Age- and time-correlated*. EM1-3
 725 have different selectivity specifications: *zero deviations*, *IID deviations*, and *AR deviations*.



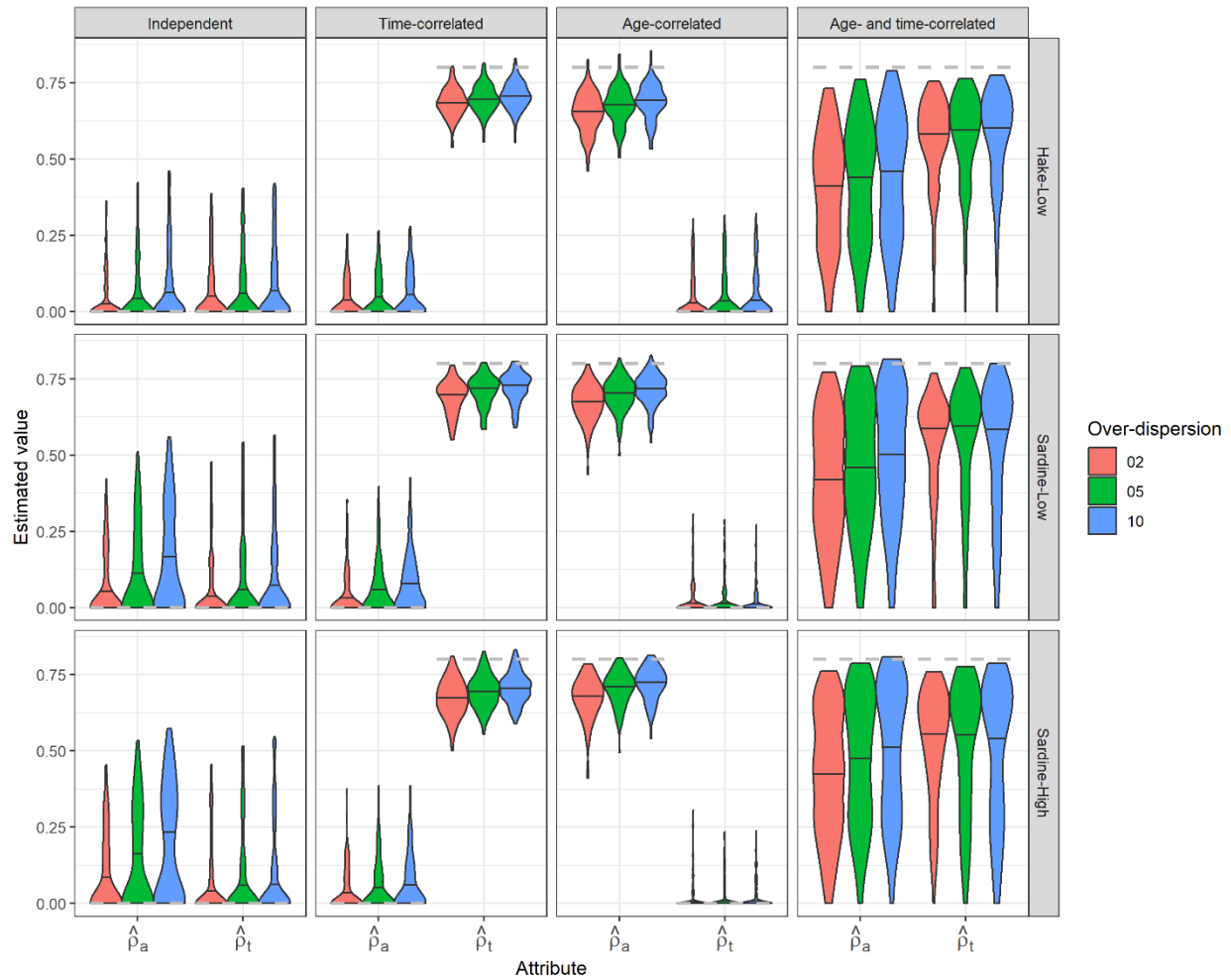
726

727 Figure 6. 3rd simulation experiment: violin plots for the ratio of effective to true sample size under
 728 three degrees of over-dispersion ($d = 2, 5, \text{ and } 10$) in age-composition data. The horizontal line in
 729 the violin plot denotes the median. The four columns correspond to the four autocorrelation cases
 730 for simulated selectivity deviations: *Independent*, *Time-correlated*, *Age-correlated*, and *Age- and*
 731 *time-correlated*. EM1-3 have different selectivity specifications: *zero deviations*, *IID deviations*,
 732 and *AR deviations*.



733

734 Figure 7. 4th simulation experiment: violin plots for the ratio of effective sample size (n_{eff}) to true
 735 sample size and the ratio of estimated ($\hat{\sigma}_S$) to the true level of selectivity variation under three
 736 degrees of over-dispersion ($d = 2, 5,$ and 10) in age-composition data. The horizontal line in the
 737 violin plot denotes the median. The four columns correspond to the four autocorrelation cases for
 738 simulated selectivity deviations: *Independent*, *Time-correlated*, *Age-correlated*, and *Age- and*
 739 *time-correlated*.



740

741 Figure 8. 4th simulation experiment: violin plots for the estimates of selectivity autocorrelations in
 742 age ($\hat{\rho}_a$) and time ($\hat{\rho}_t$) under three degrees of over-dispersion ($d = 2, 5,$ and 10) in age-composition
 743 data. The horizontal line in the violin plot denotes the median. The four columns correspond to the
 744 four autocorrelation cases for simulated selectivity deviations: *Independent*, *Time-correlated*, *Age-*
 745 *correlated*, and *Age- and time-correlated*. Horizontal dashed lines mark the true value for each
 746 autocorrelation coefficient in selectivity.

747 Table 1. Population dynamic equations in the operating model and estimation model.

No.	Equation	Comment
T1.1	$N_{a,t} = \begin{cases} R_t & a = 0 \\ N_{a-1,t-1} \exp(-S_{a,t-1}F_{t-1} - M) & 0 < a < A \\ N_{A-1,t-1} \exp(-S_{A,t-1}F_{t-1} - M) + N_{A,t-1} \exp(-S_{A,t-1}F_{t-1} - M) & a = A \end{cases}$	Stock equations
T1.2	$\ln(R_t) \sim N \left(\ln \left(\frac{4hR_0SB_t}{SB_0(1-h)+SB_t(5h-1)} \right) - \frac{\sigma_R^2}{2}, \sigma_R^2 \right)$	Recruitment
T1.3	$SB_t = \sum_{a=0}^A w_a M_a N_{a,t}$	Spawning biomass
T1.4	$C_{a,t} = N_{a,t} \frac{S_{a,t}F_t}{S_{a,t}F_t+M} (1 - e^{-S_{a,t}F_t-M})$	Catch-at-age
T1.5	$\ln(N_{a,1}) \sim N \left(\ln(R_0 e^{-aM}) - \frac{\sigma_R^2}{2}, \sigma_R^2 \right)$	Initial conditions
T1.6	$F_t = F_{t-1} \left(\frac{SB_{t-1}}{\gamma SB_0} \right)^\lambda \quad (F_{t=1} = 0.1)$	Fishing mortality

748

749 Table 2. Parameter values for the two types of life history investigated in this study.

Parameter Name	Symbol	Pacific hake	Pacific sardine
Natural mortality rate	M	0.386 yr-1	0.552 yr-1
Length at age 0	L_0	1 cm	1 cm
Asymptotic maximum length	L_{inf}	90 cm	30 cm
Von Bertalanffy growth coefficient	k	0.20 yr-1	0.30 yr-1
Log-maximum annual spawner per spawner	$LMARR$	2	1
Age at 50% selection in the fishery	S_{50}	5.44	3.55
Rate of change in selectivity at age	S_{slope}	1	2
Age at 50% maturity	a_{mat}	5.44	3.55
Steepness of the Beverton-Holt SR function	h	0.83	0.55
Ratio of equilibrium SB to unfished SB	γ	0.4	0.4
Acceleration rate in fishing mortality	λ	0.2	0.2

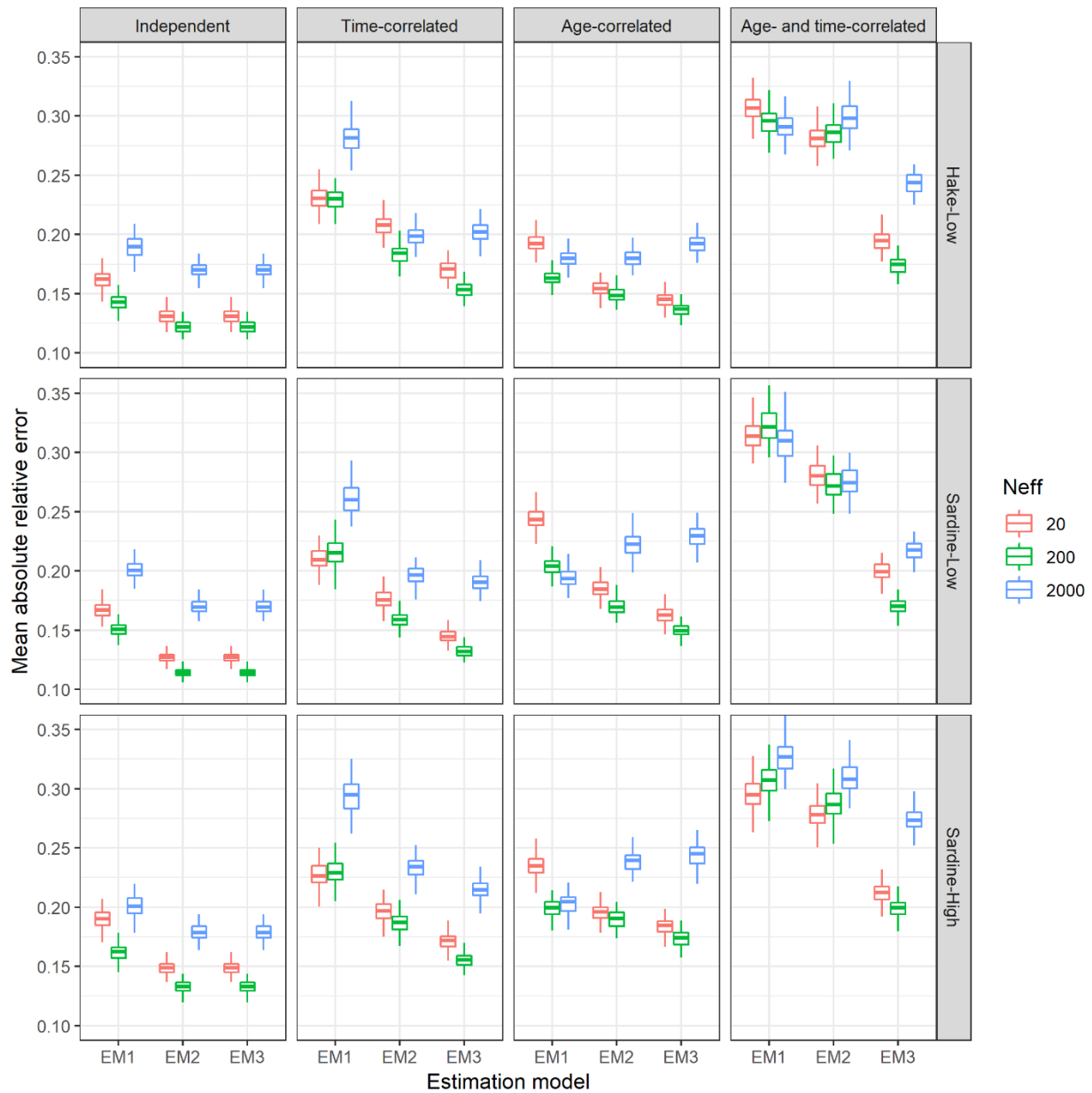
750

751 Table 3. Summary of the factorial design for each simulation experiment in this study. The
752 columns from left to right represent experiment number, operating models (1-4 represent
753 *Independent, Time-correlated, Age-correlated, and Age- and time-correlated*), estimation model
754 (1-3 represent *zero deviations, IID deviations, and AR deviations*), how the effective sample size
755 is estimated in the estimation model, data weighting methods (McAllister-Ianelli (M-I), Francis,
756 and Dirichlet-multinomial (D-M)), the degree of over-dispersion in simulated age-composition
757 data, the input sample size for the D-M method, and how the level of variation in selectivity is
758 specified in the estimation model.

Exp	OM	EM	n_{eff}	Data weighting	Over-dispersion	n_{input} (D-M)	Sel var
1	1-4	1-3	Fixed	0.1x, 1x, 10x n_{true}	-	-	true
2	1-4	1-3	Estimated	M-I, Francis, D-M	-	n_{true}	true
3	1-4	3	Estimated	D-M	2x, 5x, 10x	2x, 5x, 10x n_{true}	true
4	1-4	3	Estimated	D-M	2x, 5x, 10x	2x, 5x, 10x n_{true}	estimated

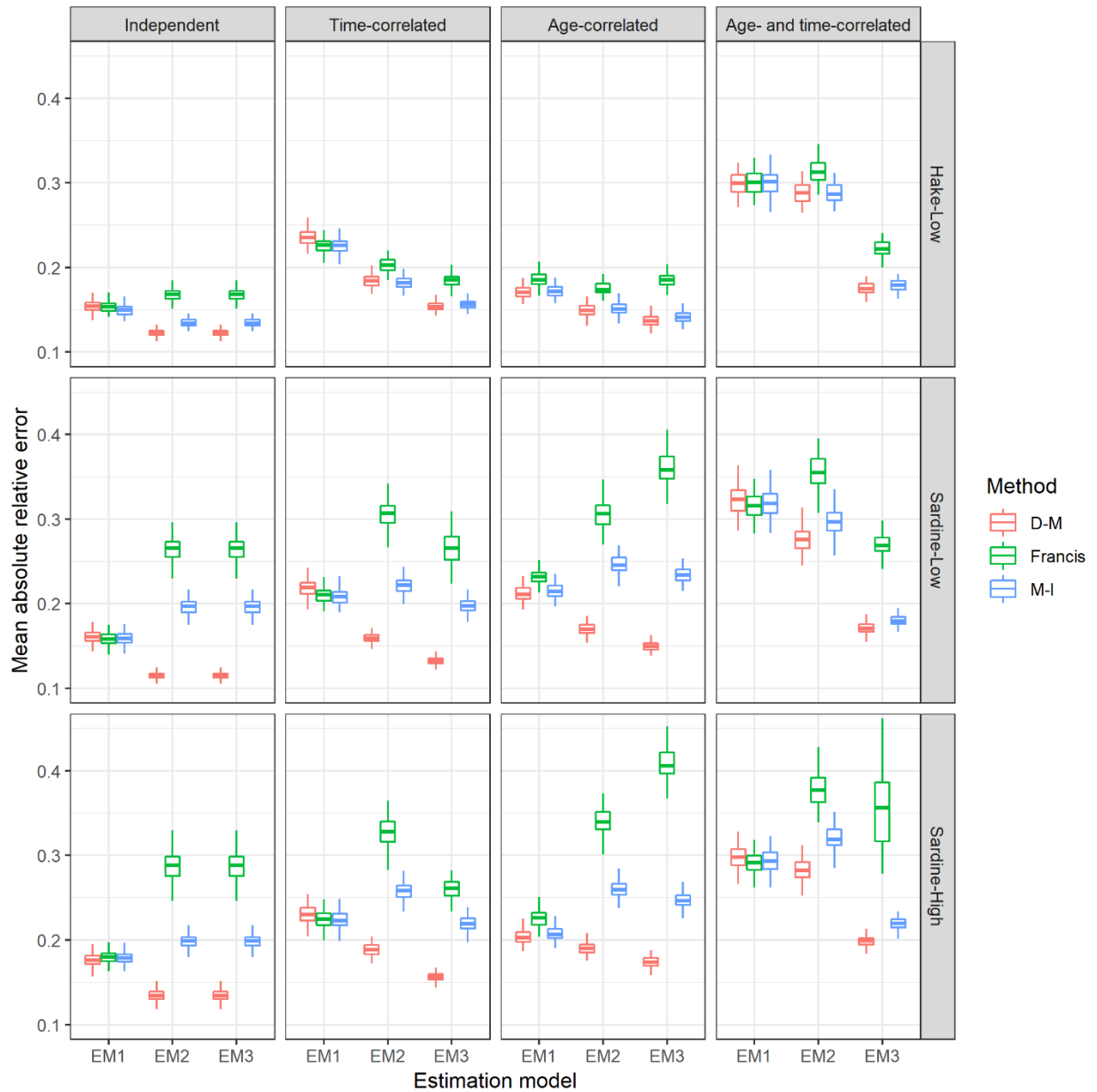
759

760



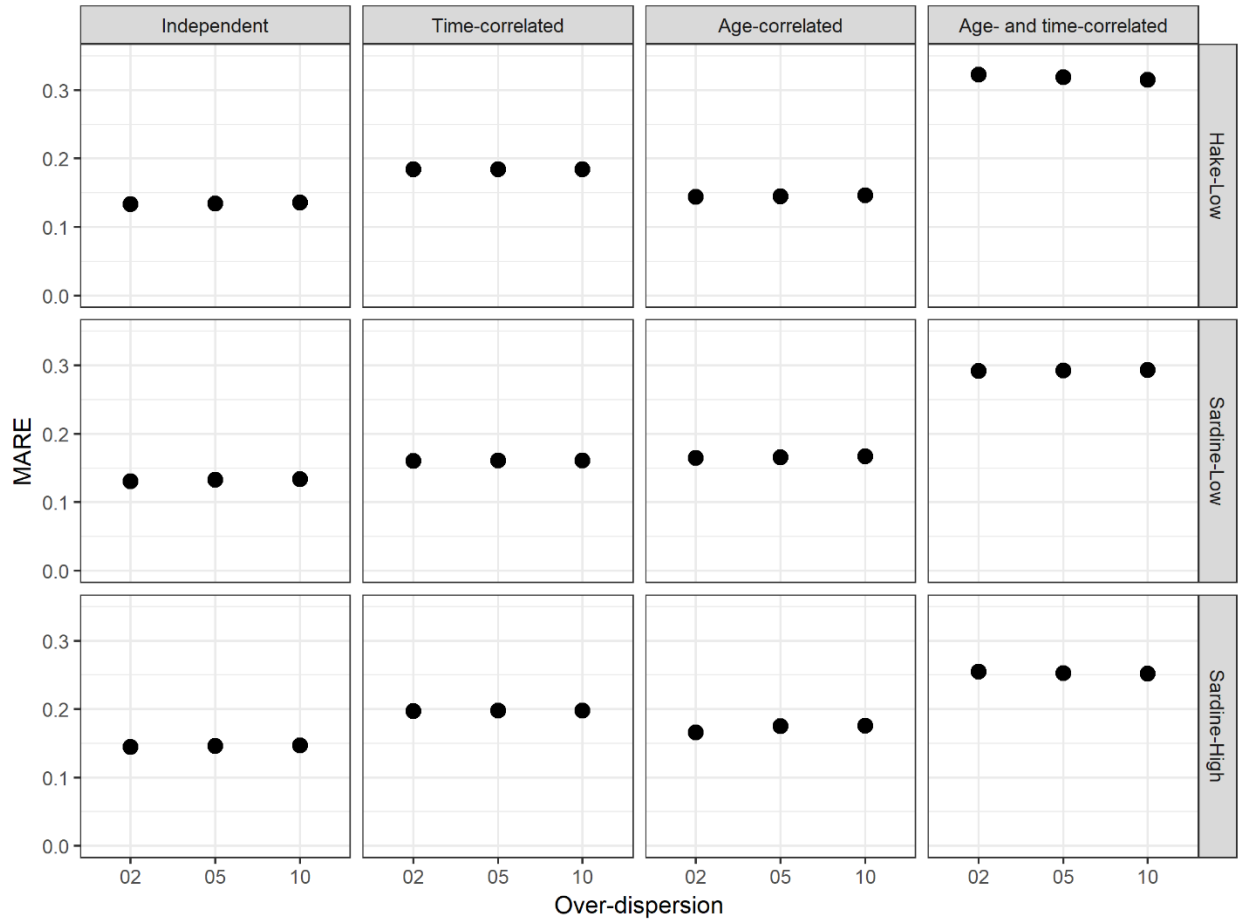
762

763 Figure A1. 1st simulation experiment: boxplot for the mean absolute relative error in the estimate
 764 of final year spawning biomass showed in Figure 3. To estimate the uncertainty of the mean
 765 absolute relative error, the 400 replicates in this simulation experiment were randomly resampled
 766 with replacement for 400 times.



767

768 Figure A2. 2nd simulation experiment: boxplot for the mean absolute relative error in the estimate
 769 of final year spawning biomass showed in Figure 5. To estimate the uncertainty of the mean
 770 absolute relative error, the 400 replicates in this simulation experiment were randomly resampled
 771 with replacement for 400 times.



772

773 Figure A3. 4th simulation experiment: mean absolute relative error in the estimate of final year
 774 spawning biomass three degrees of over-dispersion ($d = 2, 5, \text{ and } 10$) in age-composition data.

775 The four columns correspond to the four autocorrelation cases for simulated selectivity deviations:

776 *Independent, Time-correlated, Age-correlated, and Age- and time-correlated.* EM1-3 have

777 different selectivity specifications: *zero deviations, IID deviations, and AR deviations.*

778