

1 **The surprising sensitivity of index scale to delta-model**

2 **assumptions: recommendations for model-based index standardization**

3
4 James T. Thorson¹, Curry Cunningham², Elaina Jorgensen³, Andrea Havron⁴, Peter-John F.
5 Hulson⁵, Cole C. Monnahan⁶, Paul von Szalay³

6
7 1 Habitat and Ecological Processes Research Program, Resource Ecology and Fisheries
8 Management, Alaska Fisheries Science Center, National Oceanic and Atmospheric
9 Administration, 7600 Sand Point Way NE, Seattle, WA 98115, USA.

10 2 College of Fisheries and Ocean Sciences, University of Alaska Fairbanks, 17101 Point Lena
11 Loop Road, Juneau, AK 99801.

12 3 Groundfish Assessment Program, Resource Assessment and Conservation Engineering
13 Division, Alaska Fisheries Science Center, National Marine Fisheries Service, NOAA, 7600
14 Sand Point Way N.E., Seattle, WA 98115, USA.

15 4 School of Aquatic and Fishery Sciences, University of Washington, Box 355020, Seattle, WA,
16 98195, USA

17 5 Marine Ecology and Stock Assessment, Auke Bay Laboratories, Alaska Fisheries Science
18 Center, National Marine Fisheries Service, NOAA, 17109 Point Lena Loop Road, Juneau, AK
19 99801, USA

20 6 Status of Stocks and Multispecies Assessments program, Resource Ecology and Fisheries
21 Management, Alaska Fisheries Science Center, National Oceanic and Atmospheric
22 Administration, 7600 Sand Point Way NE, Seattle, WA 98115, USA.

23

24 Abstract

25 Delta-models (a.k.a. hurdle models) are widely used to fit biomass samples that include zeros
26 and a skewed response for positive catches, and spatio-temporal extensions of these models are
27 increasingly used to quantify trends in abundance (i.e., estimate abundance indices). Previous
28 research has shown estimated indices are proportional to changes in abundance. However, little
29 research has tested the performance of delta-models for estimating “scale”; that is, whether
30 abundance indices are not just proportional to population changes but also have the correct
31 absolute value. We use data for twenty species in the eastern Bering Sea and Gulf of Alaska as
32 well as a factorial experiment conditioned on data for Gulf of Alaska Pacific cod to support five
33 conclusions related to scale in spatio-temporal delta-models. First, we show that conventional
34 (nonspatial) delta-models are surprisingly sensitive to the *a priori* choice of probability
35 distribution for positive catches, where gamma and Tweedie models give similar scale estimates
36 but other distributions generally differ. Second, these same distributions also estimate widely
37 different scales when using spatio-temporal delta-models, and the delta-gamma and Tweedie
38 models provide similar scale to design-based indices. Third, model selection using marginal
39 AIC often identifies the lognormal distribution as most parsimonious, despite it resulting in
40 systematically higher abundance than design-based indices for many species. Fourth, scale is
41 sensitive to the spatial resolution (i.e., number of knots) used in fitting the spatio-temporal model
42 when using a naïve “empirical Bayes” estimator, but less sensitive when applying an epsilon
43 bias-correction estimator. Fifth, the factorial simulation experiment suggests that the Tweedie
44 and delta-gamma distributions perform well even when applied to data simulated from an
45 inverse-Gaussian or lognormal distribution, whereas the opposite is not true. We conclude that
46 index scale is sensitive to delta-model specification, and we make five recommendations when
47 using spatio-temporal delta-models for index standardization: (1) apply the epsilon or other bias-

48 correction methods to reduce sensitivity of index scale on spatio-temporal model resolution;
49 either (2) compare the scale of delta-model indices with that of design-based indices when
50 design-based indices are available or (3) use the delta-gamma or Tweedie distribution by default
51 when design-based indices are not available; (4) do not assume that AIC will identify the model
52 specification that results in the most appropriate scale; and (5) consider apparent mismatches in
53 index scale depending upon whether an assessment model specifies or estimates the associated
54 catchability coefficient and whether the design-based index is believed to measure total
55 abundance for a fully-selected age or length-class.

56

57 Keywords: Vector autoregressive spatio-temporal model; VAST; delta model; Tweedie
58 distribution; stock assessment; abundance index; catchability coefficient

59

60 1. Introduction

61 Fisheries scientists worldwide support fisheries management by estimating stock status
62 and sustainable fishing levels. They typically do this by fitting population-dynamics models to
63 fishery catches, measurements of age and length composition, and indices of population
64 abundance (Methot, 2009). Many common stock-assessment models are fitted to abundance-
65 indices that are measures (or proxies) of biomass. Fisheries scientists have therefore developed a
66 wide range of methods to sample local biomass and subsequently estimate total biomass over a
67 pre-defined spatial domain. These methods include design-based indices, which are constructed
68 from field-samples of biomass following a probabilistic design (wherein every sampling unit is
69 sampled with a pre-specified probability) and an associated statistical estimator (Cochran, 1977;
70 Smith, 1990; Petitgas, 2001). However, design-based indices are not appropriate for fishery-
71 dependent data that are not collected under a probabilistic design, or for surveys where the
72 design has changed substantially over time (i.e., adding northern stations in the eastern Bering
73 Sea, or changing the southern extent in the West Coast triennial bottom trawl survey). The
74 inability to apply design-based estimators in these instances has led to interest in model-based
75 biomass estimators, including the widely used delta-model (Pennington, 1983; Lo et al., 1992;
76 Stefansson, 1996).

77 The delta-model has been widely used for over 35 years, and separately models the
78 probability that each sample encounters a given species (termed “encounter probability” here)
79 and the probability distribution for sample biomass given that the species is encountered (termed
80 “positive catch rate” here). Aitchison (1955) originally described the delta-model as a mixture
81 distribution that contained a point mass at zero and a conditional distribution describing positive
82 (non-zero) values. The delta-lognormal distribution was proposed in the follow-up paper by
83 Aitchison and Brown (1957) and was first applied in fisheries by Pennington (1983) to describe

84 Atlantic mackerel egg production. Lambert (1992) first used a logit link function to approximate
85 the probability of encountering zero as a linear function of covariates in the context of zero-
86 inflated Poisson distributions. This approach was first applied in fisheries by Lo et al. (1992) to
87 calculate an index of relative abundance for anchovy. Delta-models were then popularized in
88 subsequent publications (Stefansson, 1996; Maunder and Punt, 2004).

89 Ongoing research has also developed spatio-temporal models that account for the
90 correlation among survey observations resulting from their proximity in space and/or time
91 (Banerjee et al., 2003; Cressie and Wikle, 2011), and these methods have recently been adapted
92 to a delta-modelling framework (Shelton et al., 2014; Thorson et al., 2015). The benefit of a
93 spatio-temporal delta-model can be seen by comparison with a design-based estimator.
94 Specifically, spatially-correlated variability in habitat quality will result in residual variance
95 among samples within each spatial stratum in a stratified-random design; this residual variance
96 will result in increasing variance for the resulting index when using a stratified-random design-
97 based estimator. In these cases, accounting for the randomized location of samples can control
98 for this spatially-correlated variability, and therefore can substantially reduce standard errors for
99 spatio-temporal indices (Shelton et al., 2014; Cao et al., 2017). In addition to increased index
100 precision, spatio-temporal delta-models have been shown to reduce biologically-implausible
101 variation in indices for long-lived species (Gertseva and Thorson 2013). When a spatio-temporal
102 delta-model was fit to U.S. West Coast trawl survey data for 28 groundfish species, confidence
103 intervals from the conventional design-based approach were 60% larger on average than those
104 derived from the spatio-temporal estimator (Thorson et al. 2015).

105 Spatio-temporal delta-models have previously and continue to be implemented for a wide
106 range of purposes. They have been used extensively for standardization of US West Coast

107 groundfish trawl survey data and are seeing increased application to Alaska groundfish survey
108 data (see list in Thorson, 2019a). In other US fisheries, spatio-temporal delta-models have been
109 implemented to estimate indices using data from multiple trawl surveys (Perretti and Thorson
110 2019), or from a mix of trawl and fixed-gear survey observations (Gruss and Thorson 2019).
111 Bayesian spatio-temporal delta-models have also been developed for standardization of
112 crustacean indices from trawl survey data from the Mediterranean Sea (Arcuti et al. 2016) and
113 shark bycatch in Canadian waters (Cosandey-Godin et al., 2014). For conservation planning,
114 spatio-temporal models have been used to integrate data from seven fisheries-independent
115 surveys, with the goal of quantifying spatial separation among target and non-target species in
116 highly-mixed Celtic Sea fisheries (Dolder et al. 2018), and to quantify spatial bycatch risk in the
117 Pacific Ocean (Stock et al. 2020). Finally, spatio-temporal delta-models have been utilized for
118 ecological inference to describe changes in species distribution, concentration, and habitat
119 association (Thorson et al. 2016a, Thorson et al. 2016b).

120 Design-based biomass indices derived from fishery-independent bottom trawl surveys are
121 fitted within many age-structured stock assessments for fish stocks in the North Pacific
122 (NPFMC, 2019a, 2019b). Age-structured models have the capacity to estimate the catchability
123 coefficient representing the ratio of predicted and index biomass (Arreguín-Sánchez, 1996).
124 Catchability coefficients are extremely influential with respect to the scale of biomass estimated
125 by a stock assessment model and are typically either estimated as a parameter or fixed at some
126 predetermined value (Wilberg et al., 2010). The estimated value for the catchability coefficient is
127 affected by spatial overlap between the stock and the spatial extent of the survey (“horizontal
128 availability”), the stocks’ vertical availability in the water column, and the stocks’ vulnerability
129 to the gear used to capture the fish (Cordue, 2007). Given the potential sensitivity of survey

130 index scale to standardization methods, and the interaction between index scale and the
131 catchability coefficient on stock assessment results, it is useful to summarize the many ways
132 catchability is specified within assessments currently.

133 We explore stock assessments at the Alaska Fisheries Science Center (AFSC) as an
134 example of stock-assessment practices for specifying the catchability coefficient throughout the
135 US and worldwide. Stock assessments at the AFSC treat the catchability coefficient using a
136 variety of approaches (see Table 1 for summary) ranging from fixing it at a value *a priori* (e.g.,
137 Bryan, 2017) to estimated freely (Thompson and Thorson, 2019). When the catchability
138 coefficient is fixed *a priori*, the survey biomass is treated as an absolute index and any change in
139 the scale of survey biomass would have direct influence on parameters that determine the scale
140 of the population (such as average recruitment and natural mortality rate). When the catchability
141 coefficient is estimated freely, the survey biomass is treated as a relative index and any
142 multiplicative change in index scale will be offset by a corresponding change in the estimated
143 catchability coefficient. Between these two extremes, some stock assessments estimate the
144 catchability coefficient using a prior distribution (either in a Bayesian or penalized likelihood
145 framework) with an associated level of uncertainty; this specified uncertainty determines the
146 degree to which the estimated catchability coefficient is able to deviate from the mean of this
147 prior distribution. When specifying a prior distribution, an infinitesimally small uncertainty is
148 equivalent to specifying a fixed value for the catchability coefficient, and an infinite level of
149 uncertainty (using a normal prior distribution with arbitrarily large variance) is equivalent to
150 freely estimating the catchability coefficient. As a consequence, the impact of changing the scale
151 of the survey index on modeled quantities from an assessment, such as spawning biomass or
152 management reference points, will be determined by the degree of precision ascribed to the

153 assumed prior on catchability coefficients: changes in index scale will be more influential on
154 modeled quantities in cases of a precise (low variance) prior and less influential in cases of
155 imprecise (high variance) prior on catchability.

156 The probability distribution for positive catches specified in a delta-model can directly
157 affect the absolute scale of the estimated index, and this is particularly important in stock
158 assessments where the catchability coefficient is fixed *a priori* or has an informative prior
159 distribution. For instance, the delta-model can result in a biased estimate of average biomass
160 when the probability distribution is mis-specified with respect to the distribution of residuals
161 (Hvingel et al., 2012; Myers and Pepin, 1990). Furthermore, delta-models can be highly
162 sensitive to deviations from model assumptions that are otherwise difficult to detect using
163 standard statistical diagnostics (Syrjala, 2000). In response, many approaches have been
164 proposed and/or applied for selecting the most appropriate distribution. Graphical tests such as
165 Taylor's power rule may help narrow the proposed set of distributions (Dick, 2004). Diagnostic
166 tests like simple Pearson correlation and normality tests on residuals, but also the lesser-known
167 Pregibon, modified Hosmer-Lemeshow, Kolmogorov-Smirnov, and Anderson-Darling tests have
168 also been explored but without consensus about their performance (Hvingel et al., 2012; Ng and
169 Cribbie, 2017). Researchers have also selected among alternative distributions using information
170 criteria like the Akaike and Bayesian Information Criteria (Akaike, 1974; Schwarz, 1978;
171 Burnham and Anderson, 2002), which appear reliable in simulations under ideal conditions and
172 sufficient sample sizes (Dick, 2004; Mitchell et al., 2015). However, sometimes AIC will select
173 models that fail diagnostic tests or can be unreliable with small sample sizes (Dick, 2004; Ng and
174 Cribbie, 2017). Furthermore, these previous simulations used GLMs without spatial effects such
175 that conclusions may not apply to spatio-temporal GLMMs. Consequently, the best statistical

176 approach for selecting the distribution for positive catch rates in spatio-temporal delta-models
177 remains unknown.

178 In this analysis, we first illustrate that the scale of an abundance-index estimated using a
179 conventional (nonspatial) delta-model is highly dependent upon the assumed distribution for
180 positive catch rates. We then compare index estimates from four spatio-temporal models (using
181 delta-gamma, delta-lognormal, delta-inverse-Gaussian, and Tweedie distributions) with design-
182 based estimates for twenty stocks in the eastern Bering Sea and Gulf of Alaska. Previous
183 research has developed an epsilon bias-correction estimator (Thorson and Kristensen, 2016) that
184 corrects for “retransformation bias” arising when random effects are transformed when
185 calculating a quantity of interest (Thorson, 2019b), but no previous study has used a simulation
186 experiment to demonstrate its importance when estimating abundance using a spatio-temporal
187 model. Similarly, we are not aware of any previous simulation study exploring how alternative
188 choices about spatial scale can affect the performance of a spatio-temporal index standardization
189 model. We therefore compare performance within a factorial design of twenty species, four
190 distributions, three spatial resolutions, and two estimators (either naïve or using the epsilon bias-
191 correction estimator). We then identify which distribution(s) provide an approximately equal
192 number of years where the abundance index is greater or less than the design-based index (i.e.
193 equivalent scale of design and model-based indices), as well as which distribution(s) estimate a
194 similar ratio between the modeled and design-based index. Finally, we use a factorial simulation
195 design conditioned upon data for Pacific cod (*Gadus macrocephalus*) in the Gulf of Alaska,
196 where we simulate data using each of the four models and fit each data set with these same four
197 estimation models. Using this simulation design, we again determine the ratio of index-scale
198 with the true population scale, as well as root-mean-squared error, to identify whether any model

200 performs best on average. Based on these findings we provide generic advice for configuring
201 delta-models for estimating abundance indices for use in stock assessments.

202 2. Methods

203 2.1 Overview

204 We seek to determine what specification for a spatio-temporal index standardization
205 model results in an index scale that matches estimates from a design-based estimator. We
206 specifically explore two alternative types of index standardization models: a delta-model
207 involving two linear predictors, or a compound Poisson-gamma (a.k.a. Tweedie) distribution
208 involving a single linear predictor. For the delta-model, we specifically explore three alternative
209 distributions for positive catch rates: a lognormal, gamma, or inverse-Gaussian distribution. This
210 then results in four model-specifications in total. All models are implemented using the Vector
211 Autoregressive Spatio-Temporal (VAST) model (Thorson and Barnett, 2017; Thorson, 2019a),
212 as implemented in package VAST release number 3.5.0 available online
213 (<https://github.com/James-Thorson-NOAA/VAST>) for the R statistical environment (R Core
214 Team, 2017). We do not explore the potential role of covariates in the following, although future
215 research could continue to explore tradeoffs associated with their inclusion (e.g., Johnson et al.,
216 2019).

217 We apply these four model specifications in two separate explorations:

- 218 1. Case study: The first is a case-study demonstration, where we fit these four model-
219 specifications to data for twenty selected species in the Gulf of Alaska and eastern Bering
220 Sea. We conduct two separate experiments using these case-study species. In the first, we fit
221 nonspatial models that estimate a separate intercept for each linear predictor in each year to
each species. This experiment is useful to show that differences in model scale arise between

222 alternative model specifications even in the simplest possible specification of an index-
223 standardization model. In the second, we fit a spatio-temporal model to data for each
224 species. In this experiment, we then compare results with a design-based estimator for each
225 species, to see which model specification results in a similar index scale to the design-based
226 estimator.

227 2. Factorial simulation experiment: The second is a factorial simulation experiment, where we
228 fit each model specification to data for a single species (Pacific cod in the Gulf of
229 Alaska). Given the estimated fixed and random effects for that species, we then simulate
230 multiple replicate data sets. For each data set, we then fit all four estimation models. This
231 then results in a 4×4 factorial cross of 4 operating models and 4 estimation models per
232 simulation replicate. We refer to scenarios where the estimation model matches the
233 operating model as a “self-test”, while other scenarios explore the implications of model mis-
234 specification on estimation model performance.

235 We describe each of these explorations in more detail below.

236 2.2 Model structure

237 In the following, we fit to observed biomass b_i for each sample i using either a Poisson-link
238 delta-model (Thorson, 2018) or a compound Poisson-gamma model (Foster and Bravington,
239 2013). Delta-models have conventionally involved a logit-linked linear predictor for encounter
240 probability, and a separate log-linked linear predictor for catch rates given an encounter
241 (Stefansson, 1996). However, we instead use a Poisson-link delta model that previous research
242 has shown to fit better while yielding a model structure that is more similar to the compound
243 Poisson-gamma distribution.

244 Poisson-link delta-models involve two log-linked linear predictors:

$$\log(n(s_i, t_i)) = \beta_n(t_i) + \omega_n^*(s_i) + \varepsilon_n^*(s_i, t_i) \quad (1)$$

$$\log(w(s_i, t_i)) = \beta_w(t_i) + \omega_w^*(s_i) + \varepsilon_w^*(s_i, t_i),$$

245 where $\beta_n(t)$ is an annually varying intercept for each modeled year $t \in \{t_{min}, \dots, t_{max}\}$, $\omega_n^*(s)$ is
 246 spatial variation that is constant over time (termed “spatial variation”) for location $s \in \Omega$ within a
 247 fixed spatial domain Ω , and ε_n^* is spatial variation that varies among years (termed “spatio-
 248 temporal variation”) in the 1st log-linked linear predictor $n(s, t)$ and similar notation is used for
 249 the second log-linked linear predictor $w(s, t)$. The product of these linear predictors $d(s, t) =$
 250 $n(s, t)w(s, t)$ is then population density $d(s, t)$ at each location s and time t . By contrast, the
 251 compound Poisson-gamma model involves a single log-linked linear predictor for density:

$$\log(d(s_i, t_i)) = \beta_d(t_i) + \omega_d^*(s_i) + \varepsilon_d^*(s_i, t_i) \quad (2)$$

252 which again includes an annual intercept, spatial, and spatio-temporal variation.

253 These models then involve specifying a probability distribution B for each sample of
 254 biomass b_i . The Poisson-linked delta-models convert $n(s_i, t_i)$ and $w(s_i, t_i)$ to encounter
 255 probability p_i and positive catch rate r_i , which varies among samples i occurring at a given
 256 location s_i and time t_i due to differences in area-swept a_i . The Poisson-linked delta-model
 257 assumes that individuals are randomly distributed in the vicinity of sampling:

$$p_i = 1 - \exp(-a_i n(s_i, t_i)) \quad (3)$$

$$r_i = \frac{a_i n(s_i, t_i) w(s_i, t_i)}{p_i}$$

258 and all delta-models assume the same probability for encounter probability:

$$\Pr(B = 0) = 1 - p_i \quad (4)$$

259 while alternative delta-models differ in the distribution for positive catches. Specifically we use
 260 a bias-corrected lognormal where dispersion parameter θ is the standard deviation in log-space:

$$\Pr(B = b_i | B > 0) = \text{Lognormal} \left(B; \log(r_i) - \frac{\theta^2}{2}, \theta^2 \right) \quad (5A)$$

261 or use a shape-scale parameterization of the Gamma distribution where dispersion θ is the
 262 coefficient of variation:

$$\Pr(B = b_i | B > 0) = \text{Gamma}(B; \theta^{-2}, r_i \theta^2) \quad (5B)$$

263 or finally we use the mean-lambda parameterization of the inverse-Gaussian distribution, where
 264 dispersion θ is again the coefficient of variation

$$\Pr(B = b_i | B > 0) = \text{Inv. Gaussian}(B; r_i, \theta^{-2}). \quad (5C)$$

265 By contrast, the compound Poisson-gamma distribution replaces Eq. 4-5 with a single
 266 distribution for biomass B

$$\Pr(B = b_i) = \text{Tweedie}(B; a_i d_i, \theta, \phi). \quad (6)$$

267 While estimating dispersion θ and power parameter $1 < \phi < 2$. Lognormal, gamma, and
 268 inverse-Gaussian distributions are all parameterized such that r_i represents the mean of positive-
 269 catch rates, such that d_i is the mean of expected catches for all distributions. However, these
 270 distributions differ somewhat in how variance is assumed to vary as a function of the mean
 271 (“mean-variance relationship”). Similarly, these distributions assign a greater or lesser
 272 probability to “extreme catches” (i.e., catches greater than ten times the expected value), and
 273 these “extreme catch events” are a well-known property of demersal fish surveys (Thorson et al.,
 274 2011). For example, the lognormal has skewness of $CV^3 + 3CV$ (where CV is the measurement
 275 error coefficient of variation) while the gamma has skewness of $2CV$. Given that the estimated
 276 CV is typically above 1.0, these distributions can have substantially different skewness. As a
 277 consequence, extremely high (or low) catches will have a greater “leverage” on predicted density
 278 for some distributions than others.

279 All models adopt a predictive-process framework for predicting spatial and spatio-
 280 temporal variation at the location s_i of each sample i , or location s_g of each extrapolation-grid
 281 cell g , given the value at n_s knots (Banerjee et al., 2008). Specifically, we specify that the value
 282 of spatial and spatio-temporal variables at each knot follows a Gaussian Markov random field:

$$\boldsymbol{\omega}_n \sim MVN(\mathbf{0}, \sigma_\omega^2 \mathbf{Q}_n^{-1}) \quad (7)$$

$$\boldsymbol{\varepsilon}_n(t) \sim MVN(\mathbf{0}, \sigma_\varepsilon^2 \mathbf{Q}_n^{-1}),$$

283 where \mathbf{Q} is a sparse precision matrix that approximates a Matern correlation function with
 284 decorrelation rate κ_n that varies among linear predictors and a transformation matrix \mathbf{H} that
 285 approximates geometric anisotropy and is shared among linear predictors. These spatial
 286 variables are then pre-multiplied by a matrix that represents bilinear interpolation (Lindgren and
 287 Rue, 2015):

$$\boldsymbol{\omega}_n^* = \mathbf{A}\boldsymbol{\omega}_n \quad (8)$$

$$\boldsymbol{\varepsilon}_n^*(t_i) = \mathbf{A}\boldsymbol{\varepsilon}_n(t_i)$$

288 and where spatial and spatio-temporal variables are treated similarly for other linear predictors
 289 $w(s, t)$ and $d(s, t)$. Specifically, interpolation matrix \mathbf{A} has a row for each extrapolation-grid
 290 cell and a column for each knot. It is nonzero for only three elements of each row (hence a
 291 “sparse” matrix), with nonzero values corresponding to the weight assigned to three vertices
 292 surrounding a given location when interpolating from three neighboring knots within a
 293 triangulated mesh.

294 Parameters are estimated by identifying the value of fixed effects that maximizes the
 295 marginal likelihood when integrated across random effects. We approximate this
 296 multidimensional integral using the Laplace approximation, as implemented using Template
 297 Model Builder (Kristensen et al., 2016). After identifying fixed effects, we then apply an

298 “empirical Bayes” estimator, which fixes random effects to their value that maximizes the joint
299 likelihood conditional on estimated fixed effects. Derived quantities can then be calculated from
300 the maximum likelihood estimate of fixed effects and empirical Bayes estimate of random
301 effects. However, derived quantities that are calculated from a nonlinear transformation of
302 random effects will be subject to “retransformation bias” when applying this naïve estimator.
303 We therefore also apply the “epsilon bias-correction estimator” that corrects for the degree of
304 nonlinearity and variance of random effects when calculating derived quantities, including
305 biomass indices (Thorson and Kristensen, 2016).

306 To estimate parameters for these models the user must:

- 307 1. Choose which probability distribution to use for the positive catches (lognormal, gamma,
308 etc.);
- 309 2. Choose the spatial resolution by specifying the number of interior knots n_x to use, which are
310 then augmented with boundary knots to determine the size of spatial and spatio-temporal
311 random effects n_s ;
- 312 3. Choose whether to use the naïve or epsilon bias-correction estimator for derived quantities.

313 We seek to provide generic guidance for these three decisions while using the “predictive
314 process” and exploring outcomes with modeled spatial resolution ranging from 100, 250, and
315 500 knots, $n_x = \{100, 250, 500\}$.

316 **2.3 Case study design**

317 Reviews for recent stock assessments at the Alaska Fisheries Science Center (AFSC) have
318 recommended further exploration of VAST regarding model specification. We therefore
319 conduct a case-study comparison of VAST models with design-based indices for twenty selected
320 species in the Gulf of Alaska and eastern Bering Sea (see Table 1 for list). The eastern Bering
321 Sea has followed a fixed-station design for bottom-trawl samples using the 83-112 gear from

322 1982-2019, where the number of samples has increased over time from approximately 350 to
323 375 per year (Lauth and Conner, 2016). The Gulf of Alaska has followed a random stratified
324 design for bottom trawl samples from 1984-2019, using the Poly Nor’eastern gear from 1990-
325 2019 and an earlier gear previously, sampling every third year from 1984-1999 and every second
326 year from 1999-2019. The number of samples per year varies from 500-850, and the sampling
327 intensity for each strata varies among years following a Neyman design based on strata-specific
328 catch rates in previous years for all species. The stratified design followed an approximately
329 consistent footprint for most years except for 2001 when the eastern Gulf of Alaska was not
330 sampled, and also in other years when deep-water strata were dropped due to funding limitations
331 (von Szalay and Raring, 2016).

332 For each of these stocks, we first fit model-based estimators that include only the annual
333 intercept in each year (β) and exclude the spatial and spatio-temporal terms (ω and ε), resulting
334 in a simple unstratified delta-model. We do not expect this specification of model-based indices
335 to accurately measure population biomass because this specification ignores spatial stratification
336 and other concerns about sampling design. However, we compare model-based indices for
337 alternative models to demonstrate the extent to which index scale can differ even when fitting a
338 simple index model.

339 For each stock, we next extract a design-based estimator using standard protocols and
340 software for these two regions (Wakabayashi et al., 1985). We compare these with spatio-
341 temporal model-based estimators that extrapolate density to the “standard” footprint of these
342 surveys. The spatio-temporal estimator specifically predicts density at the centroid of grid cells
343 within a 2km by 2km square extrapolation-grid that serves as “quadrature points” for integrating
344 across density. This includes 36,140 grid cells for the “Eastern Bering Sea” extrapolation-grid

345 and 23,339 grid cells for the “Gulf of Alaska” extrapolation-grid; each is included in package
346 VAST and was developed previously by Angie Grieg (personal communication; retired from
347 Alaska Fisheries Science Center). We expect that these spatio-temporal model-based estimators
348 will appropriately account for spatial variation in inclusion probability (i.e., due to stratified
349 sampling) given that this probability-sampling design is constructed based on results for a wide
350 variety of species and is likely to be independent of density for any single species (Conn et al.,
351 2017).

352 The design-based estimator will be an unbiased estimator for the portion of population
353 biomass that is available to the survey in each year. We acknowledge that the design-based
354 estimator will in many cases not be an accurate representation of fully-selected abundance or
355 biomass, e.g., in cases when the stock moves into and out of the spatial footprint of a single
356 survey (Ianelli et al., 2019), moves vertically out of the area accessible to bottom trawls
357 (Kotwicki et al., 2015), or moves into areas where gear performs poorly (Thorson et al., 2013).
358 Previous studies have evaluated performance for spatio-temporal models via comparison to
359 stock-assessment model output (e.g., Cao et al., 2017; Thorson and Haltuch, 2018), but have not
360 used a simulation experiment to compare performance against the scale of design-based indices.
361 We therefore evaluate the model performance for estimating population scale relative to design-
362 based indices by we calculating the average across n_t years for both the design-based index $\bar{B} =$
363 $\frac{1}{n_t} \sum_{t=1}^{n_t} B_t$ and each model-based index $\bar{I} = \frac{1}{n_t} \sum_{t=1}^{n_t} I_t$. We then calculate the ratio of these two
364 averages $R = \bar{I}/\bar{B}$ and record this ratio for each species and model specification, for each model
365 resolution and when using either the naïve or epsilon-bias correction estimator. We seek to
366 determine what model specification results in a similar scale to design-based indices and

367 therefore identify a well performing model as one with a ratio R that is evenly distributed around
368 one, indicating that the scale is similar on average to the scale of the design-based index.

369 For each model and spatial resolution, we also calculate the Akaike Information Criterion
370 (AIC), calculated using the Laplace approximation to the marginal likelihood and the number of
371 fixed effects. We specifically seek to determine whether AIC consistently favors any model
372 specification, and if the model specification selected using this criteria varies with changes in
373 spatial resolution.

374 **2.4 Design for factorial simulation experiment**

375 We also explore model performance by conducting a 4×4 factorial design of all four model
376 specifications as both operating model and estimation model, when fixed and random effects for
377 each operating model are determined by fitting them to the bottom trawl survey data for Pacific
378 cod in the Gulf of Alaska. We note that the epsilon bias-correction estimator is computationally
379 expensive using the predictive-process model formulation, and therefore facilitate parameter
380 estimation within the replicated design by decreasing the number of extrapolation-grid cells. We
381 specifically use a k-means algorithm to identify 2000 locations, and calculate their area as the
382 sum of areas for those extrapolation-grid cells that are nearest to each. This procedure therefore
383 integrates across density using 2000 “quadrature points” rather than the original 36,140
384 extrapolation-grid cells. This decreases the spatial resolution used when integrating density, and
385 substantially reduces computation time in particular during the epsilon bias-correction estimator.
386 By using this new technique for both the estimation and operating model, we decrease the time
387 required for each simulation replicate by approximately 75%, and exploratory testing confirms
388 that it does not introduce any bias when applied to both estimation and operating models.

389 We evaluate model performance by recording the true biomass \tilde{B}_{mrt} in each operating
390 model m , simulation replicate r , and year t , and comparing this true biomass with the estimated

391 biomass $I_{mrt,d}$ for each each replicate, operating model, year, and estimation model d . We
392 specifically calculate relative error $E_{m,r,t,d} = (I_{m,r,t,d} - \tilde{B}_{m,r,t})/\tilde{B}_{m,r,t}$ and then visualize the
393 average relative error across all years and replicates for a given operating and estimation model.
394 A well-performing model will have a relative error centered on zero and a low root-mean-
395 squared relative error. In particular a minimax estimator suggests that the best model is that
396 which minimizes the maximum error across all model scenarios (Lehmann and Casella, 1998 pg.
397 309), in this case constituted by the four operating models.

398 We also evaluate model performance by calculating the correlation between the natural
399 logarithm of true density (from the operating model) and predicted density (from the estimation
400 model). In particular, we calculate the correlation separately for each year, and then average
401 across years for a given simulation replicate; this calculation emphasizes model performance in
402 identifying areas with high or low density. This comparison specifically addresses whether a
403 particular estimation model performs better or worse at identifying spatial variation in density;
404 we speculate that a different estimation model might be appropriate for accurately estimating
405 spatial variation vs. estimating the scale when integrating across space for calculating an
406 abundance index.

407 3. Results

408 Applying a nonspatial delta-model to biomass samples for twenty species in the Gulf of
409 Alaska and eastern Bering Sea shows many cases where model specification has large effects on
410 resulting index variability and scale (Fig. 1). For example, *Sebastes polyspinus* in the Gulf of
411 Alaska shows an approximately stable index using the lognormal delta-model and an increasing
412 trend for the inverse-Gaussian. By contrast, both gamma and Tweedie models show large spikes
413 in estimated abundance in 2001 and 2013, and agree with the lower abundance in 2015-2019

414 estimated by the lognormal distribution rather than the elevated estimates of the inverse-
415 Gaussian. Similarly, *S. alutus* in the Gulf of Alaska and both *Lepidopsetta polyxystra* and
416 *Limanda aspera* in the eastern Bering Sea show similar indices for gamma and Tweedie models,
417 but differ from indices arising from either lognormal or inverse-Gaussian distributions. These
418 and other examples show that sensitivity to the assumed distribution of positive catch rates is a
419 general characteristic of delta-models, rather than an issue specifically with spatio-temporal
420 delta-models.

421 Next we compare spatio-temporal indices using three resolutions (100, 250, or 500 knots)
422 with design-based indices. Illustrating results for three selected species shows that gamma and
423 Tweedie models generate similar indices, which are also similar in terms of both variability and
424 trend to the design-based indices (Fig. 2). However, models with lower resolutions (100 knots)
425 tend to estimate a higher scale than increased resolutions (250 or 500 knots) or the design-based
426 indices. For these species, the inverse-Gaussian and lognormal models produce indices that
427 show similar index trends and variability to other models and design-based indices, but differ
428 greatly in terms of scale as a function of the specified spatial resolution.

429 Notably, AIC selects the lognormal and inverse-Gaussian for 8-11 of the twenty species
430 for these three resolutions (Fig. 3), and often selects the lognormal even for species where the
431 Tweedie and gamma result in indices that have an index scale more similar to design-based
432 indices (e.g., *Sebastes alutus* in the Gulf of Alaska in Fig. 2). Specifically, the ratio of average
433 biomass for model-and design-based indices is 0.98 and 1.01 when using bias-correction and
434 high resolution for the gamma and Tweedie models, while this ratio is 1.23 and 1.60 for the
435 lognormal and inverse-Gaussian models (Fig. 4, black numbers in right column). The difference
436 between design- and model-based scale increases for the gamma and Tweedie models either

437 without epsilon bias-correction (e.g., red values in Fig. 4), or with decreasing resolution (e.g., left
438 and middle columns in Fig. 4).

439 Finally, the factorial simulation design confirms that models generally have good
440 performance (i.e., small bias and low root-mean-squared error) when the simulation and
441 estimation model have matching specification (i.e., diagonal panels in Fig. 5). However, the
442 estimation models (Fig. 5 columns) differ greatly in terms of average performance when applied
443 to data from a mis-specified simulation model. For example, the inverse-Gaussian estimation
444 model has poor performance (e.g., large positive bias) when applied to data simulated using a
445 gamma or Tweedie distribution, and the lognormal distribution also shows a smaller but still
446 substantial positive bias for these operating models. By contrast, the gamma and Tweedie
447 estimation models have a bias between -4 to +1% when applied to data for any of the operating
448 models. We therefore conclude that both gamma and Tweedie estimation models are identified
449 by a “minimax” estimator as the estimation models that minimizes the maximum error across
450 alternative operating models. By contrast, the lognormal estimation model performs somewhat
451 better than the gamma and Tweedie models with respect to the correlation between true and
452 estimated density, particularly when fitted to data generated by an inverse-Gaussian distribution
453 (Fig. 6). However, we note that all three distributions all do well in general as estimation models
454 (correlation > 0.84 for each operating model). We therefore conclude that the optimal
455 distribution for estimating spatial variation in density will in some cases be different than the
456 optimal distribution for estimating the scale of an abundance index that is in agreement with a
457 design-based estimator.

458 4. Discussion

459 In this study, we have shown that delta-gamma and Tweedie distributions result in a
460 similar scale for model-based abundance indices as design-based indices for twenty stocks in the
461 North Pacific. Results also highlight that index scale is sensitive to the number of knots used to
462 approximate spatial variation within a spatio-temporal model when using a naïve estimator, but
463 this sensitivity is mitigated when using the epsilon bias-correction estimator that accounts for
464 retransformation bias. Using the highest resolution and bias-correction estimator, the delta-
465 gamma and Tweedie models have an average ratio of 0.98 and 1.01 relative to design-based
466 indices, indicating that they have a similar scale on average to a design-based estimator. When
467 averaging design and model-based indices across years, the root-mean-squared log-ratio between
468 these averages is 0.16 and 0.24, respectively. This suggests that the difference in scale (i.e.,
469 difference in average value for design- and model-based indices for a given species) is
470 approximately 20% between these alternative approaches. Similarly, a factorial simulation
471 design suggests that delta-gamma and Tweedie models have minimal error even for data
472 simulated using other distributions, and therefore minimize the maximum error arising from
473 these candidate forms of model mis-specification. This result is similar to classical statistical
474 studies aimed at comparing lognormal and gamma distributions within generalized linear models
475 in general (Firth, 1988; Wiens, 1999). Finally, the lognormal distribution performs best
476 (followed closely by gamma and Tweedie models) at estimating spatial variation in density,
477 indicating that difficulties in estimating index scale are largely separate from model ability to
478 accurately identify spatial variation in density.

479 Spatio-temporal models fitted to biomass samples are already seeing widespread use in
480 stock, ecosystem, habitat, and climate-vulnerability assessments (Thorson, 2019a). In particular,

481 model-based indices can be generated using data that do not strictly follow a probabilistic design
482 (Ye and Dennis, 2009), or can account for failures to consistently implement a planned design.
483 However, there is more to learn regarding the expected performance of delta-models when the
484 estimation model is mis-specified with respect to the data-generating process. In particular, we
485 are surprised by the strong dependence of abundance-index scale upon the choice of probability
486 distribution for positive catch rates. Previous simulation studies have not highlighted this model
487 sensitivity because they: (1) focused on the proportionality of index estimates and true
488 abundance and thereby ignored scale (Dick, 2004; Thorson et al., 2015); (2) eliminated model
489 mis-specification by using the same distribution for generating and estimation (Johnson et al.,
490 2019); (3) explored bias for a single class of delta-model without comparing performance across
491 distributions (Myers and Pepin, 1990; Smith, 1990); (4) focused simulation testing on features
492 other than the process used to generate data used in index standardization (Berg et al., 2014; Lo
493 et al., 1992); or (5) did not document this mismatch in scale even when the estimation and
494 simulation models were mismatched (Ono et al., 2015). We recommend further testing of delta-
495 models using a variety of operating models, including individual- and agent-based models whose
496 properties will not exactly match any simple estimation model. Using a variety of operating
497 models will allow a more complete picture of the magnitude of errors arising from mis-
498 specifying the distribution for positive catch rates. We also recommend further exploration of
499 optimal ways of generating the SPDE mesh used in INLA and VAST; we have not explored this
500 in detail here, but it could be one line of research to explore the sensitivity of index scale to the
501 specified resolution.

502 The appropriate use of information criteria such as AIC in hierarchical (e.g., spatio-
503 temporal) models is an unresolved topic in statistics due to the difficulty in estimating the

504 effective degrees of freedom associated with random-effects that are shrunk towards zero
505 (Hodges and Sargent, 2001; Wikle et al., 2019 Chapter 6). Marginal AIC is defined as the AIC
506 score when counting only fixed effects, while conditional AIC is defined as AIC while partially
507 counting random effects based on their estimated variance (Vaida and Blanchard, 2005). Both
508 marginal and conditional AIC have known types of poor behavior for mixed-effects models
509 (Grevin and Kneib, 2010), and our results confirm poor behavior for marginal AIC, which
510 tended to select the lognormal distribution even in cases when its scale differed greatly from a
511 design-based estimate. Multiple methods have also been proposed to improve performance for
512 marginal and conditional AIC (Müller et al., 2013; Watanabe, 2013). For example, Shang and
513 Cavanaugh (2008) developed a bootstrap method to calculate a more appropriate penalty term,
514 Sakamoto (2019) developed a computationally efficient approach to correct for issues in the
515 marginal AIC, and Grevin and Kneib (2010) developed an analytic correction to the conditional
516 AIC.

517 In addition to model selection, new GLMM methods can allow for more rigorous model
518 validation through hypothesis testing. The DHARMA R package (Hartig, 2017) offers a suite of
519 tests and validation diagnostics to evaluate uniform residuals calculated from the empirical
520 distribution function of simulated values for an observation evaluated at the observation value.
521 One-step-ahead residuals are calculated iteratively by evaluating marginal likelihoods of
522 observation subsets against predicted values (Thygesen et al., 2017). Residuals can be compared
523 with specified distributions using tests such as the Shapiro-Wilk, Komogoroc-Smirnov or
524 Anderson-Darling hypothesis tests. However, we recommend further research regarding
525 quantitative tools for model selection and validation, to automate the process of identifying an
526 appropriate distribution for positive catch rates in spatio-temporal delta-models.

527 We also recommend continued research to identify delta-model specifications that are
528 less sensitive to likelihood choice. One idea is to develop and implement new generalized
529 distributions in VAST that contain common distributions as nested submodels, thereby replacing
530 a (categorical) model selection with (continuous) parameter estimation. Hvingel et al. (2012)
531 used the generalized gamma distribution, which adds a third parameter to the gamma and
532 contains the lognormal, gamma, Weibull, and exponential distributions as special cases (Stacy,
533 1962). This distribution is difficult to fit because its parameters are highly correlated (Stacy and
534 Mihram, 1965), although there has also been some success with reparameterizations (Prentice,
535 1974). An alternative approach would be to use robust estimators that are designed to be
536 insensitive to data drawn from a range of distributions (Maronna et al., 2019). Conceptually, a
537 robust delta-lognormal estimator would minimize sensitivity to outliers, thereby serving as a
538 reliable default. Some theoretical and simulation work has shown promise for models without
539 covariates or other effects like space (Rosales, 2009), but research is needed to extend robust
540 estimators to mixed-effects models like VAST. We encourage future studies to investigate these
541 ideas as potential solutions to make estimation of absolute indices more stable and reliable.

542 Whether to use a model- or design-based survey index in a given stock assessment
543 depends in part upon how the resulting index is subsequently treated within the assessment
544 model. In particular, it depends upon whether the index is viewed as absolute (i.e., the
545 catchability coefficient is fixed a priori), or if the survey index is treated as relative and the
546 parameter(s) describing survey catchability are estimated. Differences in index trend between
547 model and design-based indices would be important regardless of how the catchability
548 coefficient is treated, but large differences in trend were not observed among estimation spatio-
549 temporal delta-model specifications explored (e.g., Fig. 2). Differences in index scale between

550 model and design-based indices are important if the assessment treats the index as absolute, but
551 have limited impact on model results if the catchability coefficient is freely estimated. In
552 practice, bottom trawl survey biomass indices at the AFSC typically fall somewhere on a
553 continuum between absolute (q fixed at 1) and relative (q freely estimated) indices, with several
554 assessments residing somewhere in between by specifying informative priors or likelihood
555 penalties for q (Table 1). Delta-models using a gamma or Tweedie distribution generally differ
556 from the design-based index scale by 10%, and this is usually within the standard deviation of
557 the prior distribution assumed in Alaskan groundfish assessment models implementing an
558 informative prior for q (Table 1).

559 Based upon our results and in light of issues noted above, we recommend the following
560 practices when using spatio-temporal delta-models to generate abundance indices for use in stock
561 assessments:

- 562 1. *Compare model-based index scale with design-based indices when possible:* Most
563 importantly, our simulation and case-study examples highlight that the choice of distribution
564 for positive catch rates can have large effect on estimated scale. In most cases, we envision
565 that analysts will trust the scale from a design-based estimator, and that similarity in scale
566 could be one criterion (among others) for selecting among potential distributions.
- 567 2. *Use the gamma or Tweedie distributions by default when it is not possible to compare with*
568 *design-based scale:* In other cases, a design-based estimator may not be feasible, either
569 because the data are opportunistic (i.e., fishery-dependent catches), the survey substantially
570 departed from the planned design (i.e., a vessel broke down), or the design is not sufficient
571 for inference about a given stock (i.e., data from multiple designs must be combined). In

572 these cases, our simulation experiment suggests that the gamma or Tweedie distribution have
573 reasonable performance across a range of data-generating mechanisms.

574 3. *Correct for retransformation bias using the epsilon estimator:* Our case-study results suggest
575 that the epsilon bias-correction estimator (Thorson and Kristensen, 2016) results in a much
576 better match between model- and design-based index scale than the naïve empirical Bayes
577 estimator, and decreases sensitivity to model resolution.

578 4. *Do not assume that AIC is the only criterion for model performance:* Our results also suggest
579 that AIC will select the lognormal distribution even in cases where it has poor match to the
580 scale of the design-based index. We therefore recommend multiple considerations (including
581 index scale and diagnostics) when selecting a model. We also recommend future research to
582 develop automated approaches to calculate conditional AIC for models implemented in
583 Template Model Builder, including the VAST model used here. This development would
584 then allow for a detailed performance comparison between marginal and AIC for index-
585 standardization models.

586 5. *Consider assessment-model structure when deciding between model- and design-based*
587 *indices:* Finally, we note a variety of practices for treating the catchability coefficient for
588 stock assessments in the North Pacific, and suspect that this same variation arises in other
589 management regions. Eight of the twenty case-study species use a catchability coefficient
590 that is fixed a priori, and these assessments are likely to be highly sensitive to differences in
591 index scale. In cases where a design-based index is available and believed to measure total
592 abundance/biomass for a fully-selected age/length class (i.e., not missing entire spatial strata
593 due to operational problems or gear restrictions), we encourage analysts to compare the scale
594 of model-based indices with that of design-based indices and use this information to inform

595 their choice of which method to use. Six assessments estimate the catchability coefficient
596 freely, and index scale will have no effect for these assessments; in these cases, comparison
597 of scale between model- and design-based indices could be used as a diagnostic of the spatio-
598 temporal model, but will have direct impact on assessment-model results. Finally, six are
599 estimated with a prior or penalty, and prior/penalty standard deviation is typically larger than
600 the expected difference in scale between model- and design-based indices for gamma and
601 Tweedie distributions. In summary, we recommend that the index scale be compared
602 between model- and design-based indices in all three cases. However, the match in scale is
603 most important for assessments that assume a fixed catchability coefficient, and is relevant to
604 consider in cases where the design-based index is believed to measure total
605 abundance/biomass for a fully-selected age or length-class. We recognize that this
606 recommendation requires contextual information to interpret, and recommend further
607 research regarding situations when a model-based index is likely to provide a more useful
608 estimate of scale (whether due to improved precision, accounting for densities in areas that
609 are not measured within a design-based estimator, or other reasons).

610 Finally, we continue to recommend that regional authorities for scientific review establish
611 regional “Terms of Reference” (Thorson, 2019a) such that criteria for model specification are
612 clear, transparent, and easily replicated for any stock assessment within a given region.

613 **5. Acknowledgements**

614 We thank Kasper Kristensen, Hans Skaug, and the TMB development team, without which
615 VAST would not be computationally feasible. We also thank the many scientists and volunteers
616 who have contributed to the bottom trawl surveys in the Gulf of Alaska and eastern Bering Sea

617 shelf. Finally, we thank Lewis Barnett, C. O’Leary, and two anonymous reviewers for helpful
618 comments on an earlier draft.

619 6. References

- 620 Aitchison, J., 1955. On the Distribution of a Positive Random Variable Having a Discrete Probability Mass
621 at the Origin. *J. Am. Stat. Assoc.* 50, 901. <https://doi.org/10.2307/2281175>
- 622 Aitchison, J., Brown, J.A., 1957. The lognormal distribution with special reference to its uses in
623 economics. Cambridge University Press, Cambridge, MA.
- 624 Akaike, H., 1974. New look at statistical-model identification. *IEEE Trans. Autom. Control* AC19, 716–723.
- 625 Arreguín-Sánchez, F., 1996. Catchability: a key parameter for fish stock assessment. *Rev. Fish Biol. Fish.*
626 6, 221–242. <https://doi.org/10.1007/BF00182344>
- 627 Banerjee, S., Carlin, B.P., Gelfand, A.E., 2003. Hierarchical modeling and analysis for spatial data, 1st ed.
628 Chapman & Hall/CRC, Boca Raton, FL.
- 629 Banerjee, S., Gelfand, A.E., Finley, A.O., Sang, H., 2008. Gaussian predictive process models for large
630 spatial data sets. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 70, 825–848.
631 <https://doi.org/10.1111/j.1467-9868.2008.00663.x>
- 632 Berg, C.W., Nielsen, A., Kristensen, K., 2014. Evaluation of alternative age-based methods for estimating
633 relative abundance from survey data in relation to assessment models. *Fish. Res.* 151, 91–99.
634 <https://doi.org/10.1016/j.fishres.2013.10.005>
- 635 Bryan, M.D., 2017. Assessment of the northern and southern rock sole (*Lepidopsetta polyxstra* and
636 *bilineata*) stocks in the Gulf of Alaska (NPFMC Bering Sea and Aleutian Islands SAFE). North
637 Pacific Fishery Management Council, Anchorage, AK.
- 638 Burnham, K.P., Anderson, D., 2002. Model Selection and Multi-Model Inference, 2nd ed. Springer, New
639 York.
- 640 Cao, J., Thorson, J.T., Richards, R.A., Chen, Y., 2017. Spatiotemporal index standardization improves the
641 stock assessment of northern shrimp in the Gulf of Maine. *Can. J. Fish. Aquat. Sci.* 74, 1781–
642 1793. <https://doi.org/10.1139/cjfas-2016-0137>
- 643 Cochran, W.G., 1977. Sampling Techniques, 3rd Edition, 3rd ed. John Wiley & Sons.
- 644 Conn, P.B., Thorson, J.T., Johnson, D.S., 2017. Confronting preferential sampling when analysing
645 population distributions: diagnosis and model-based triage. *Methods Ecol. Evol.* 8, 1535–1546.
646 <https://doi.org/10.1111/2041-210X.12803>
- 647 Cordue, P.L., 2007. A note on non-random error structure in trawl survey abundance indices. *ICES J.*
648 *Mar. Sci.* 64, 1333–1337. <https://doi.org/10.1093/icesjms/fsm134>
- 649 Cosandey-Godin, A., Krainski, E.T., Worm, B., Flemming, J.M., 2014. Applying Bayesian spatiotemporal
650 models to fisheries bycatch in the Canadian Arctic. *Can. J. Fish. Aquat. Sci.* 72, 186–197.
651 <https://doi.org/10.1139/cjfas-2014-0159>
- 652 Cressie, N., Wikle, C.K., 2011. Statistics for spatio-temporal data. John Wiley & Sons, Hoboken, New
653 Jersey.
- 654 Dick, E.J., 2004. Beyond “lognormal versus gamma”: discrimination among error distributions for
655 generalized linear models. *Fish. Res.* 70, 351–366. <https://doi.org/10.1016/j.fishres.2004.08.013>
- 656 Firth, D., 1988. Multiplicative Errors: Log-Normal or Gamma? *J. R. Stat. Soc. Ser. B Methodol.* 50, 266–
657 268. <https://doi.org/10.1111/j.2517-6161.1988.tb01725.x>
- 658 Foster, S.D., Bravington, M.V., 2013. A Poisson–Gamma model for analysis of ecological non-negative
659 continuous data. *Environ. Ecol. Stat.* 20, 533–552. <https://doi.org/10.1007/s10651-012-0233-0>

660 Greven, S., Kneib, T., 2010. On the behaviour of marginal and conditional AIC in linear mixed models.
661 Biometrika 97, 773–789. <https://doi.org/10.1093/biomet/asq042>

662 Hartig, F., 2017. DHARMA: residual diagnostics for hierarchical (multi-level/mixed) regression models. R
663 Package Version 01 5.

664 Hodges, J.S., Sargent, D.J., 2001. Counting degrees of freedom in hierarchical and other richly-
665 parameterised models. Biometrika 88, 367–379. <https://doi.org/10.1093/biomet/88.2.367>

666 Hvingel, C., Kingsley, M.C.S., Sundet, J.H., 2012. Survey estimates of king crab (*Paralithodes*
667 *camtschaticus*) abundance off northern Norway using GLMs within a mixed generalized gamma-
668 binomial model and Bayesian inference. ICES J. Mar. Sci. 69, 1416–1426.
669 <https://doi.org/10.1093/icesjms/fss116>

670 Ianelli, J.N., Fissel, B., Holsman, K., Honkalehto, T., Kotwicki, S., Monnahan, C., Siddon, E., Stienessen, S.,
671 Thorson, J.T., 2019. Assessment of the walleye pollock stock in the Eastern Bering Sea (NPFMC
672 Bering Sea and Aleutian Islands SAFE). North Pacific Fishery Management Council, Anchorage,
673 AK.

674 Johnson, K.F., Thorson, J.T., Punt, A.E., 2019. Investigating the value of including depth during
675 spatiotemporal index standardization. Fish. Res. 216, 126–137.
676 <https://doi.org/10.1016/j.fishres.2019.04.004>

677 Kotwicki, S., Horne, J.K., Punt, A.E., Ianelli, J.N., 2015. Factors affecting the availability of walleye pollock
678 to acoustic and bottom trawl survey gear. ICES J. Mar. Sci. J. Cons. 72, 1425–1439.

679 Kristensen, K., Nielsen, A., Berg, C.W., Skaug, H., Bell, B.M., 2016. TMB: Automatic differentiation and
680 Laplace approximation. J. Stat. Softw. 70, 1–21. <https://doi.org/10.18637/jss.v070.i05>

681 Lambert, D., 1992. Zero-Inflated Poisson Regression, With an Application to Defects in Manufacturing.
682 Technometrics 34, 1–14. <https://doi.org/10.1080/00401706.1992.10485228>

683 Lauth, R.R., Conner, J., 2016. Results of the 2013 eastern Bering Sea continental shelf bottom trawl
684 survey of groundfish and invertebrate resources (NOAA Technical Memorandum No. NMFS-
685 AFSC-331). Alaska Fisheries Science Center, Seattle, WA.

686 Lehmann, E.L., Casella, G., 1998. Theory of Point Estimation, 2nd edition. ed. Springer, New York.

687 Lindgren, F., Rue, H., 2015. Bayesian spatial modelling with r-inla. J. Stat. Softw. 63, 1–25.
688 <https://doi.org/10.18637/jss.v063.i19>

689 Lo, N.C., Jacobson, L.D., Squire, J.L., 1992. Indices of Relative Abundance from Fish Spotter Data based
690 on Delta-Lognormal Models. Can. J. Fish. Aquat. Sci. 49, 2515–2526.

691 Maronna, R.A., Martin, R.D., Yohai, V.J., Salibián-Barrera, M., 2019. Robust Statistics: Theory and
692 Methods, 2 edition. ed. Wiley, Hoboken, NJ.

693 Maunder, M.N., Punt, A.E., 2004. Standardizing catch and effort data: a review of recent approaches.
694 Fish. Res. 70, 141–159. <https://doi.org/10.1016/j.fishres.2004.08.002>

695 Methot, R.D., 2009. Stock Assessment: Operational Models in Support of Fisheries Management, in:
696 Beamish, R.J., Rothschild, B.J. (Eds.), The Future of Fisheries Science in North America. Springer
697 Netherlands, Dordrecht, pp. 137–165.

698 Mitchell, E.M., Lyles, R.H., Schisterman, E.F., 2015. Positing, fitting, and selecting regression models for
699 pooled biomarker data. Stat. Med. 34, 2544–2558. <https://doi.org/10.1002/sim.6496>

700 Müller, S., Scealy, J.L., Welsh, A.H., 2013. Model Selection in Linear Mixed Models. Stat. Sci. 28, 135–167.
701 <https://doi.org/10.1214/12-STS410>

702 Myers, R.A., Pepin, P., 1990. The robustness of lognormal-based estimators of abundance. Biometrics
703 46, 1185–1192.

704 Ng, V.K.Y., Cribbie, R.A., 2017. Using the Gamma Generalized Linear Model for Modeling Continuous,
705 Skewed and Heteroscedastic Outcomes in Psychology. Curr. Psychol. 36, 225–235.
706 <https://doi.org/10.1007/s12144-015-9404-0>

707 NPFMC, 2019a. Stock Assessment and Fishery Evaluation Report for the Groundfish Resources of the
708 Gulf of Alaska. North Pacific Fishery Management Council, Anchorage, AK.

709 NPFMC, 2019b. Stock Assessment and Fishery Evaluation Report for the Groundfish Resources of the
710 Bering Sea and Aleutian Islands Region. North Pacific Fishery Management Council, Anchorage,
711 AK.

712 Ono, K., Punt, A.E., Hilborn, R., 2015. Think outside the grids: An objective approach to define spatial
713 strata for catch and effort analysis. *Fish. Res.* 170, 89–101.
714 <https://doi.org/10.1016/j.fishres.2015.05.021>

715 Pennington, M., 1983. Efficient Estimators of Abundance, for Fish and Plankton Surveys. *Biometrics* 39,
716 281–286.

717 Petitgas, P., 2001. Geostatistics in fisheries survey design and stock assessment: models, variances and
718 applications. *Fish Fish.* 2, 231–249.

719 Prentice, R.L., 1974. A Log Gamma Model and Its Maximum Likelihood Estimation. *Biometrika* 61, 539–
720 544. <https://doi.org/10.2307/2334737>

721 R Core Team, 2017. R: A Language and Environment for Statistical Computing. R Foundation for
722 Statistical Computing, Vienna, Austria.

723 Rosales, M.A.C., 2009. The Robustness of Confidence Intervals for the Mean of Delta Distribution.
724 Western Michigan University.

725 Sakamoto, W., 2019. Bias-reduced marginal Akaike information criteria based on a Monte Carlo method
726 for linear mixed-effects models. *Scand. J. Stat.* 46, 87–115.

727 Schwarz, G., 1978. Estimating the Dimension of a Model. *Ann. Stat.* 6, 461–464.
728 <https://doi.org/10.1214/aos/1176344136>

729 Shang, J., Cavanaugh, J.E., 2008. Bootstrap variants of the Akaike information criterion for mixed model
730 selection. *Comput. Stat. Data Anal.* 52, 2004–2021. <https://doi.org/10.1016/j.csda.2007.06.019>

731 Shelton, A.O., Thorson, J.T., Ward, E.J., Feist, B.E., 2014. Spatial semiparametric models improve
732 estimates of species abundance and distribution. *Can. J. Fish. Aquat. Sci.* 71, 1655–1666.
733 <https://doi.org/10.1139/cjfas-2013-0508>

734 Smith, S.J., 1990. Use of statistical models for the estimation of abundance from groundfish trawl survey
735 data. *Can. J. Fish. Aquat. Sci.* 47, 894–903.

736 Stacy, E.W., 1962. A Generalization of the Gamma Distribution. *Ann. Math. Stat.* 33, 1187–1192.
737 <https://doi.org/10.1214/aoms/1177704481>

738 Stacy, E.W., Mihram, G.A., 1965. Parameter Estimation for a Generalized Gamma Distribution.
739 *Technometrics* 7, 349–358. <https://doi.org/10.1080/00401706.1965.10490268>

740 Stefansson, G., 1996. Analysis of groundfish survey abundance data: combining the GLM and delta
741 approaches. *ICES J Mar Sci* 53, 577–588.

742 Syrjala, S.E., 2000. Critique on the use of the delta distribution for the analysis of trawl survey data. *ICES*
743 *J. Mar. Sci.* 57, 831–842. <https://doi.org/10.1006/jmsc.2000.0571>

744 Thompson, G., Thorson, J.T., 2019. Assessment of the Pacific cod stock in the Eastern Bering Sea. In
745 Stock assessment and fishery evaluation report for the groundfish resources of the Bering Sea
746 and Aleutian Islands (NPFMC Bering Sea and Aleutian Islands SAFE). North Pacific Fishery
747 Management Council, Anchorage, AK.

748 Thorson, J.T., 2019a. Guidance for decisions using the Vector Autoregressive Spatio-Temporal (VAST)
749 package in stock, ecosystem, habitat and climate assessments. *Fish. Res.* 210, 143–161.
750 <https://doi.org/10.1016/j.fishres.2018.10.013>

751 Thorson, J.T., 2019b. Perspective: Let’s simplify stock assessment by replacing tuning algorithms with
752 statistics. *Fish. Res., Recruitment: Theory, Estimation, and Application in Fishery Stock*
753 *Assessment Models* 217, 133–139. <https://doi.org/10.1016/j.fishres.2018.02.005>

754 Thorson, J.T., 2018. Three problems with the conventional delta-model for biomass sampling data, and a
755 computationally efficient alternative. *Can. J. Fish. Aquat. Sci.* 75, 1369–1382.
756 <https://doi.org/10.1139/cjfas-2017-0266>

757 Thorson, J.T., Barnett, L.A.K., 2017. Comparing estimates of abundance trends and distribution shifts
758 using single- and multispecies models of fishes and biogenic habitat. *ICES J. Mar. Sci.* 74, 1311–
759 1321. <https://doi.org/10.1093/icesjms/fsw193>

760 Thorson, J.T., Haltuch, M.A., 2018. Spatiotemporal analysis of compositional data: increased precision
761 and improved workflow using model-based inputs to stock assessment. *Can. J. Fish. Aquat. Sci.*
762 1–14. <https://doi.org/10.1139/cjfas-2018-0015>

763 Thorson, J.T., Kristensen, K., 2016. Implementing a generic method for bias correction in statistical
764 models using random effects, with spatial and population dynamics examples. *Fish. Res.* 175,
765 66–74. <https://doi.org/10.1016/j.fishres.2015.11.016>

766 Thorson, J.T., M. Elizabeth, C., Stewart, I.J., Punt, A.E., 2013. The implications of spatially varying
767 catchability on bottom trawl surveys of fish abundance: a proposed solution involving
768 underwater vehicles. *Can. J. Fish. Aquat. Sci.* 70, 294–306.

769 Thorson, J.T., Shelton, A.O., Ward, E.J., Skaug, H.J., 2015. Geostatistical delta-generalized linear mixed
770 models improve precision for estimated abundance indices for West Coast groundfishes. *ICES J.*
771 *Mar. Sci. J. Cons.* 72, 1297–1310. <https://doi.org/10.1093/icesjms/fsu243>

772 Thorson, J.T., Stewart, I.J., Punt, A.E., 2011. Accounting for fish shoals in single-and multi-species survey
773 data using mixture distribution models. *Can. J. Fish. Aquat. Sci.* 68, 1681–1693.

774 Thygesen, U.H., Albertsen, C.M., Berg, C.W., Kristensen, K., Nielsen, A., 2017. Validation of ecological
775 state space models using the Laplace approximation. *Environ. Ecol. Stat.* 24, 317–339.
776 <https://doi.org/10.1007/s10651-017-0372-4>

777 Vaida, F., Blanchard, S., 2005. Conditional Akaike information for mixed-effects models. *Biometrika* 92,
778 351–370.

779 von Szalay, P.G., Raring, N.W., 2016. Data report: 2015 Gulf of Alaska bottom trawl survey (NOAA
780 Technical Memorandum No. NMFS-AFSC-325). US Department of Commerce, National Oceanic
781 and Atmospheric Administration, National Marine Fisheries Service, Alaska Fisheries Science
782 Center, Seattle, WA.

783 Wakabayashi, K., Bakkala, R.G., Alton, M.S., 1985. Methods of the U.S.-Japan demersal trawl surveys, in:
784 Bakkala, R.G., Wakabayashi, K. (Eds.), *Results of Cooperative US-Japan Groundfish Investigations*
785 *in the Bering Sea during May-August 1979*.

786 Watanabe, S., 2013. A widely applicable Bayesian information criterion. *J. Mach. Learn. Res.* 14, 867–
787 897.

788 Wiens, B.L., 1999. When Log-Normal and Gamma Models Give Different Results: A Case Study. *Am. Stat.*
789 53, 89–93. <https://doi.org/10.1080/00031305.1999.10474437>

790 Wikle, C.K., Zammit-Mangion, A., Cressie, N., 2019. *Spatio-Temporal Statistics with R*, 1 edition. ed.
791 Chapman and Hall/CRC, Boca Raton.

792 Wilberg, M.J., Thorson, J.T., Linton, B.C., Berkson, J., 2010. Incorporating time-varying catchability into
793 population dynamic stock assessment models. *Rev. Fish. Sci.* 18, 7–24.

794 Ye, Y., Dennis, D., 2009. How reliable are the abundance indices derived from commercial catch-effort
795 standardization? *Can. J. Fish. Aquat. Sci.* 66, 1169–1178.

796

797 **Figures and Tables**

798

799 Table 1: All stocks included in analysis, including the scientific and common name of the assessed species, the region for each stock
800 (GOA=Gulf of Alaska, EBS=Eastern Bering Sea), and a reference for the stock assessment. We also list how the catchability
801 coefficient for the bottom trawl survey is treated (either fixed at a value *a priori*, estimated with a prior distribution, or estimated
802 freely without a prior distribution), the coefficient of variation for the associated prior when estimated using one, and whether
803 catchability is varying over time either through a time-dependent parameterization or implicit variation due to estimated time-varying
804 selectivity.

Scientific name	Common name	Region	Assessment reference	Treatment of catchability coefficient	CV of prior on catchability coefficient	Time-varying catchability
<i>Atheresthes stomias</i>	Arrowtooth Flounder	GOA	Spies et al., 2019a	Fixed	--	Not time-dependent
<i>Microstomus pacificus</i>	Dover Sole	GOA	McGilliard et al., 2019	Fixed and estimated with prior	85%	Time-blocks (fixed one block, estimated one block)
<i>Hippoglossoides elassodon</i>	Flathead Sole	GOA	Turnock et al., 2017	Fixed	--	Not time-dependent
<i>Sebastes polyspinis</i>	Northern Rockfish	GOA	Cunningham et al., 2018	Estimated with prior	45%	Not time-dependent
<i>Gadus macrocephalus</i>	Pacific Cod	GOA	Barbeaux et al., 2019	Estimated freely	--	Time-dependent through selectivity
<i>Sebastes alutus</i>	Pacific Ocean Perch	GOA	Hulson et al., 2019	Estimated with prior	45%	Not time-dependent

<i>Lepidopsetta polyxystra</i> and <i>L. bilineata</i>	Northern and Southern Rock Sole	GOA	Bryan, 2017	Fixed	--	Not time-dependent
<i>Gadus chalcogrammus</i>	Walleye Pollock	GOA	Dorn et al., 2019	Estimated with prior	10%	Not time-dependent
<i>Pleuronectes quadrituberculatus</i>	Alaska Plaice	EBS	Wilderbuer and Nichol, 2019	Fixed	--	Not time-dependent
<i>Beringraja binoculata</i>	Alaska Skate	EBS	Ormseth, 2018	Fixed	--	Not time-dependent
<i>Atheresthes stomias</i>	Arrowtooth Flounder	EBS	Spies et al., 2019a	Estimated freely	--	Time-dependent through annual deviations related to bottom water temperature
<i>Reinhardtius hippoglossoides</i>	Greenland Turbot	EBS	Bryan et al., 2018a	Fixed	--	Not time-dependent
<i>Atheresthes evermanni</i>	Kamchatka Flounder	EBS	Bryan et al., 2018b	Estimated freely	--	Time-dependent through annual deviations related to bottom water temperature
<i>Lepidopsetta polyxystra</i>	Northern Rock Sole	EBS	Wilderbuer et al., 2018	Fixed	--	Not time-dependent
<i>Gadus macrocephalus</i>	Pacific Cod	EBS	Thompson and Thorson, 2019	Estimated freely	--	Time-dependent through selectivity
<i>Hippoglossus stenolepis</i>	Pacific Halibut	EBS	--	Estimated freely in areas-as-fleets model	--	Not time-dependent
<i>Gadus chalcogrammus</i>	Walleye Pollock	EBS	Ianelli et al., 2019	Estimated freely	--	Time-dependent through selectivity
<i>Limanda aspera</i>	Yellowfin Sole	EBS	Spies et al., 2019b	Estimated with prior	90%	Time-dependent through annual deviations related to bottom water temperature

<i>Anoplopoma fimbria</i>	Sablefish	GOA and EBS	Hanselman et al., 2019	Estimated with prior	30%	Not time-dependent
---------------------------	-----------	----------------	---------------------------	-------------------------	-----	--------------------

805

806

807 **Figure captions**

808

809 Figure 1: Model-based abundance indices (y-axis) in each year (x-axis) for each of twenty
810 species (panels), showing estimates from four nonspatial models: three Poisson-link delta-
811 models using lognormal (red), gamma (green), and inverse-Gaussian (blue) distributions for
812 positive catches, and a Tweedie distribution for modeling both encounter rate and positive catch
813 rate (grey).

814

815 Figure 2: Visualizing model-based abundance indices (y-axis, shown on log-scale) in each year
816 (x-axis) for each of three species (columns) using four alternative distributions (rows), where
817 each panel shows the abundance index (line) and 95% confidence interval (shaded area) for three
818 different spatial resolutions (see color legend in bottom-right panel indicating the number of
819 knots) as well as the design-based estimators (black dots), and each panel also includes the
820 percent AIC weight for each distribution and resolution across models (e.g., where percentages
821 for a given color sum to 100% for each column)

822

823 Figure 3: Marginal AIC weights (y-axis) for each distribution (x-axis) using a given model
824 resolution (rows). Each bar includes multiple colored segments, showing the AIC weight for
825 each individual stock.

826

827 Figure 4: Histogram showing number of species (y-axis) with a given ratio between model- and
828 design-based indices when each is averaged across years (x-axis, shown on log-scale) for three
829 model resolutions (columns) and distributions (rows). A well-performing model will have an

830 average ratio near 0 on the log scale or 1.0 on the linear scale. Each panel also has a set of
831 numbers showing the average ratio (top-left, where 1.0 corresponds to a similar scale) and the
832 root-mean-squared error (top-right, where 0.0 corresponds to a scale that is identical between
833 model- and design-based approaches) when using epsilon bias-correction (black) or not using
834 bias-correction (red).

835

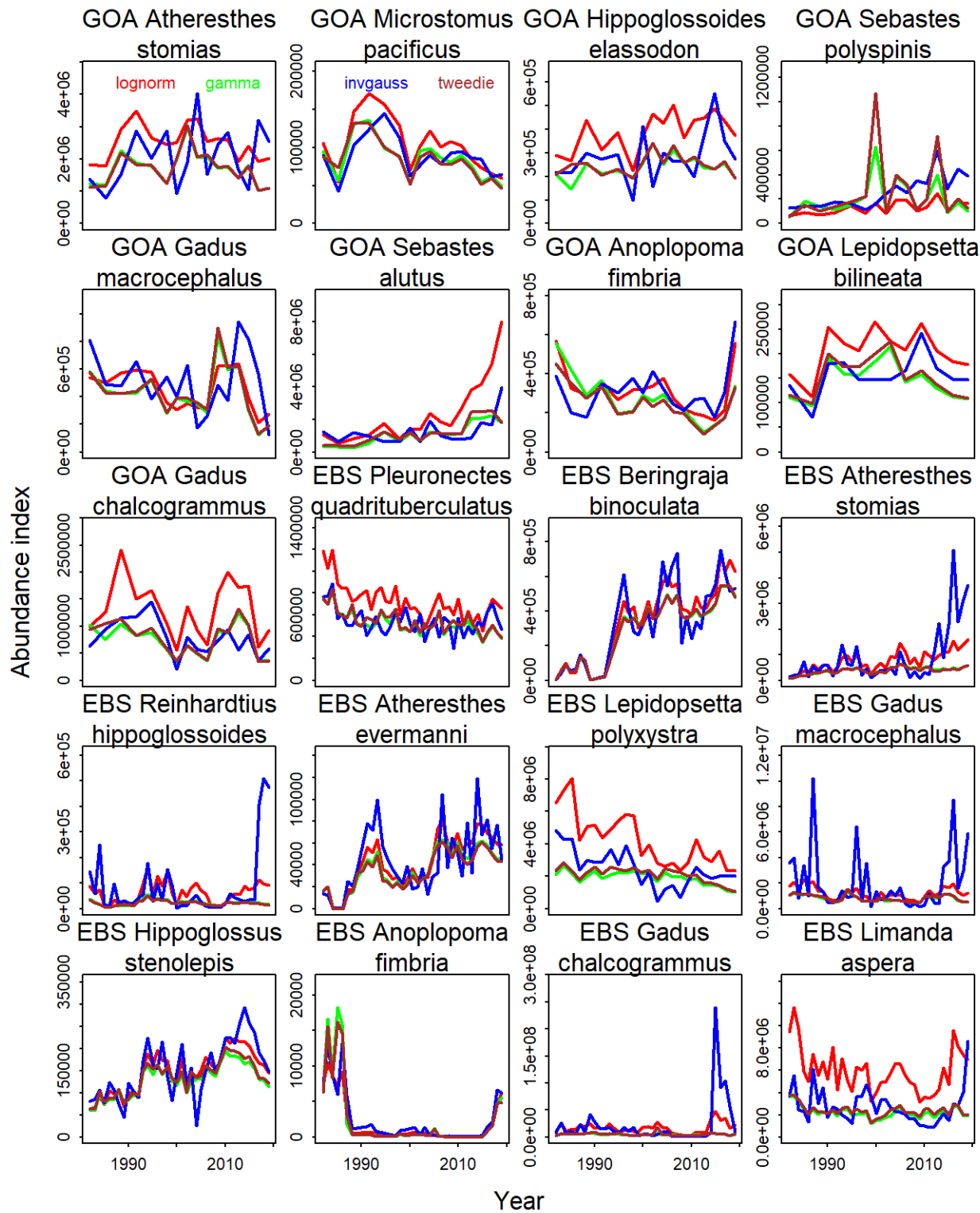
836 Figure 5: Distributions of relative errors when comparing estimated and true abundance indices
837 (x-axis) within a factorial simulation experiment conditioned on survey data for Pacific cod in
838 the Gulf of Alaska, where the four distributions are used as operating models (rows, such that
839 they are fitted to available data where fixed and random effects are then held constant when
840 simulating new sampling data following the same sampling design), as well as estimation models
841 (columns, i.e., fitted to simulated data from a given operating model). Panels on the diagonal
842 involve the same estimation and operating model and are expected to have low error, while each
843 column shows the performance of a given estimation model across different forms of model mis-
844 specification. A generally well-performing estimation model will have a relative error near 0
845 (dashed vertical line) for all panels in a given column; each panel also lists the bias and root-
846 mean-square-error (in parentheses) calculated for all replicates for a given operating and
847 estimation model.

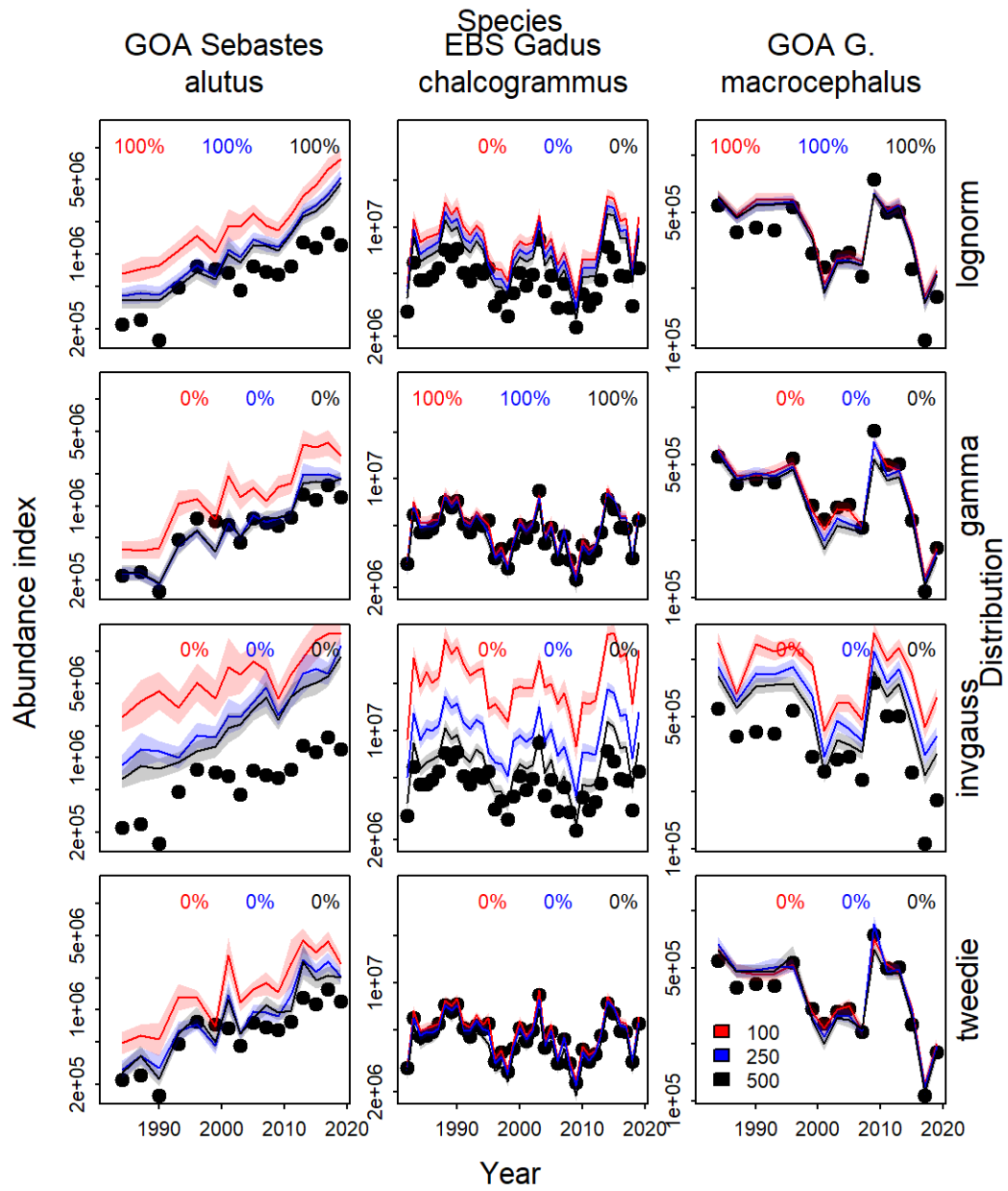
848

849 Figure 6: Distribution of Pearson correlation coefficients between estimated and true density,
850 calculated for each year individually and then averaged across years for a given simulation
851 replicate (x-axis), where the four distributions are used as operating models (rows) as well as
852 estimation models (columns). See Fig. 5 caption for more details. A well-performing estimation

853 model will have a correlation near 1.0 for each panels in a given column; each panel also lists the
854 average correlation calculated for all replicates for a given operating and estimation model.

855

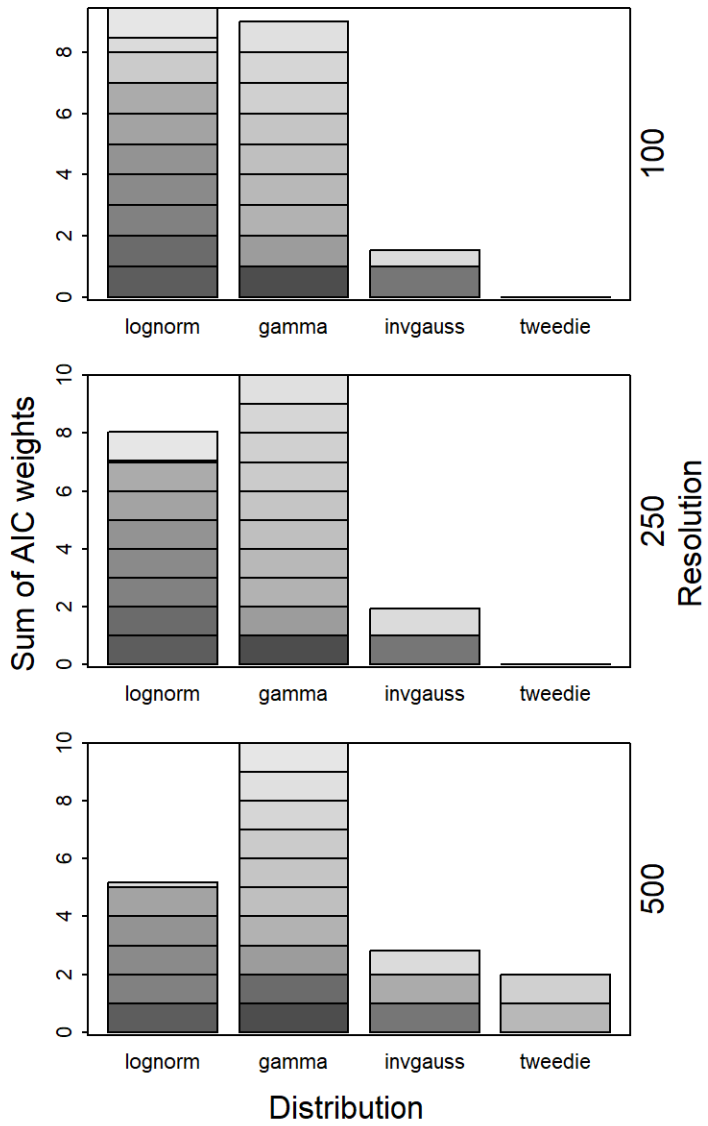




860

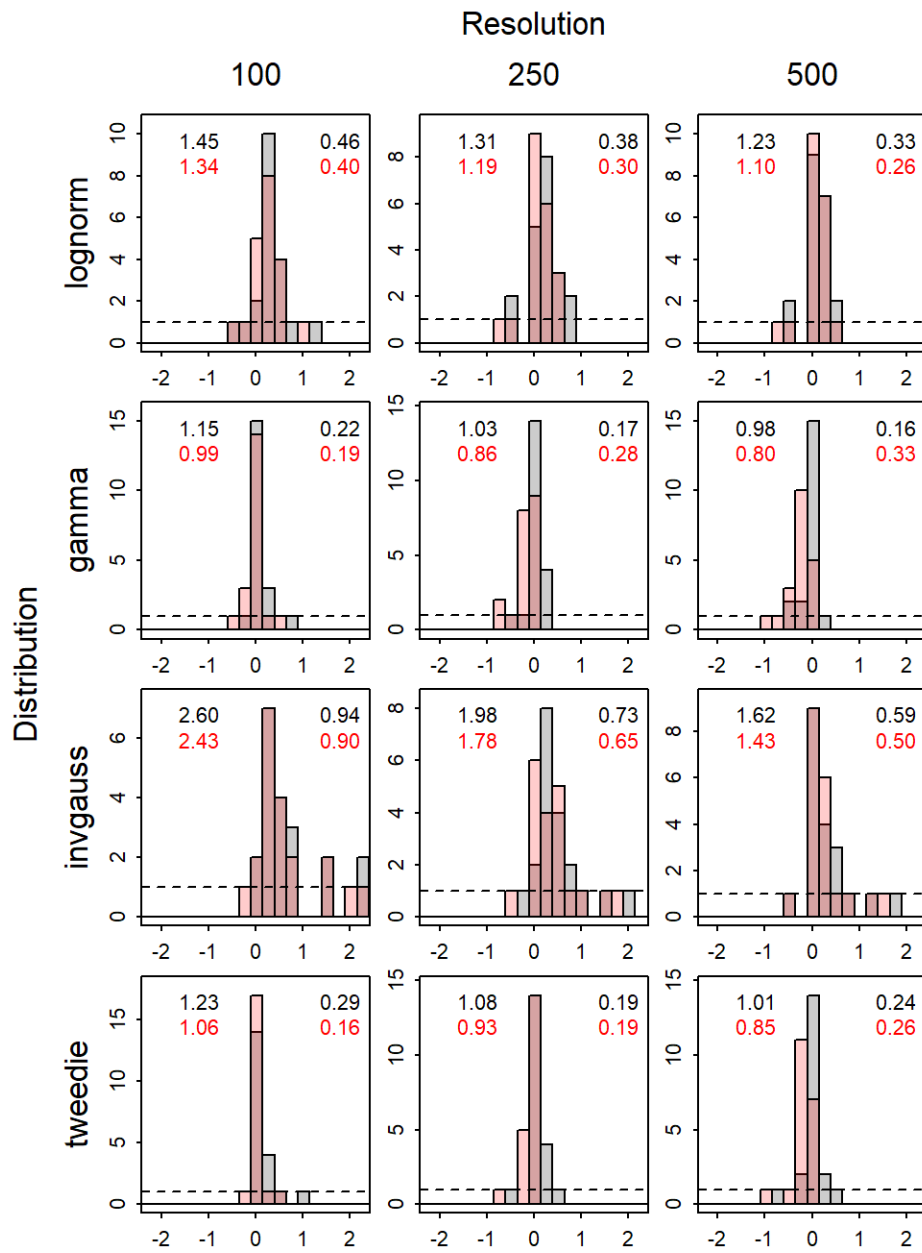
861

862 Fig. 3



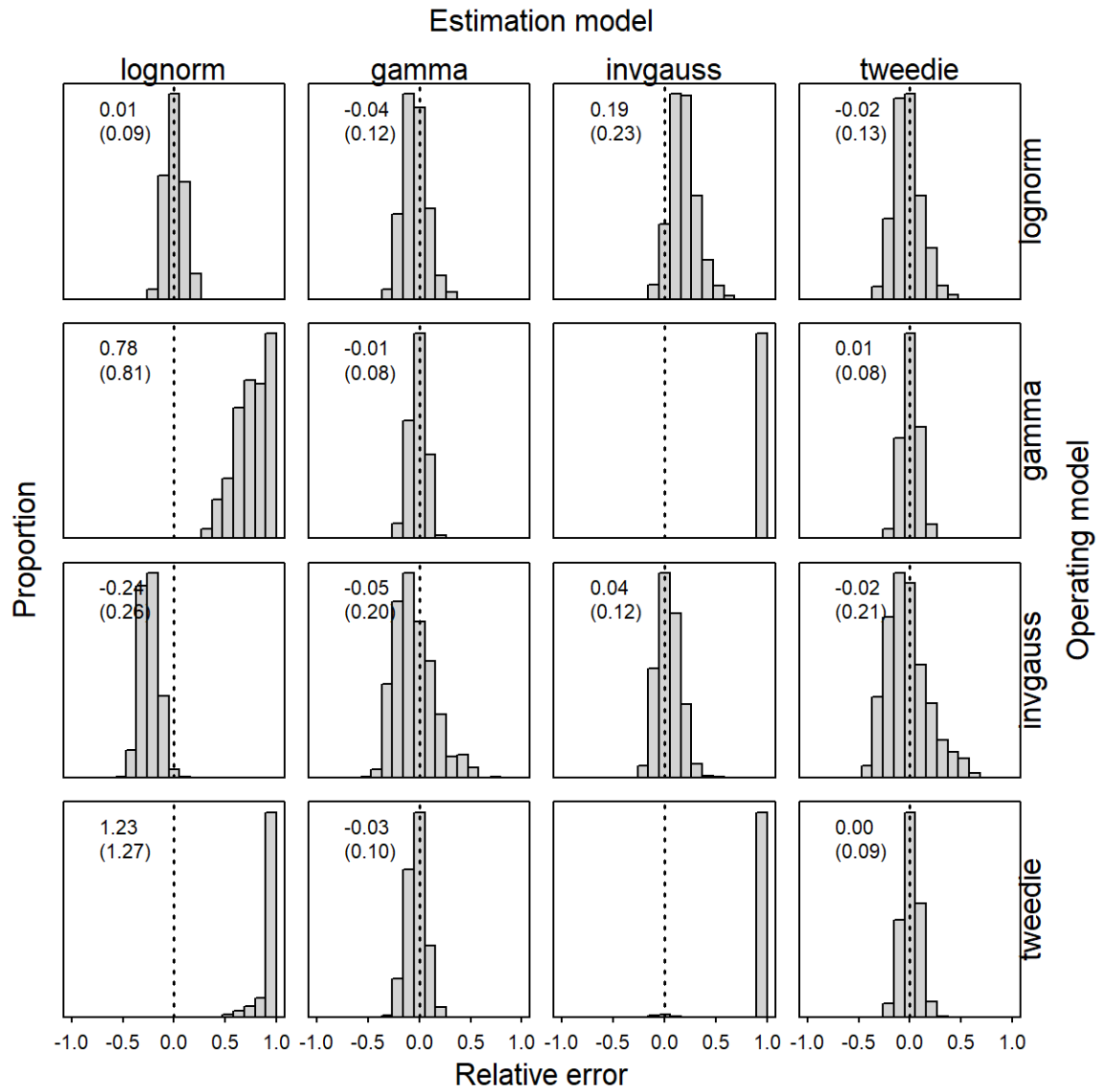
863

864



866

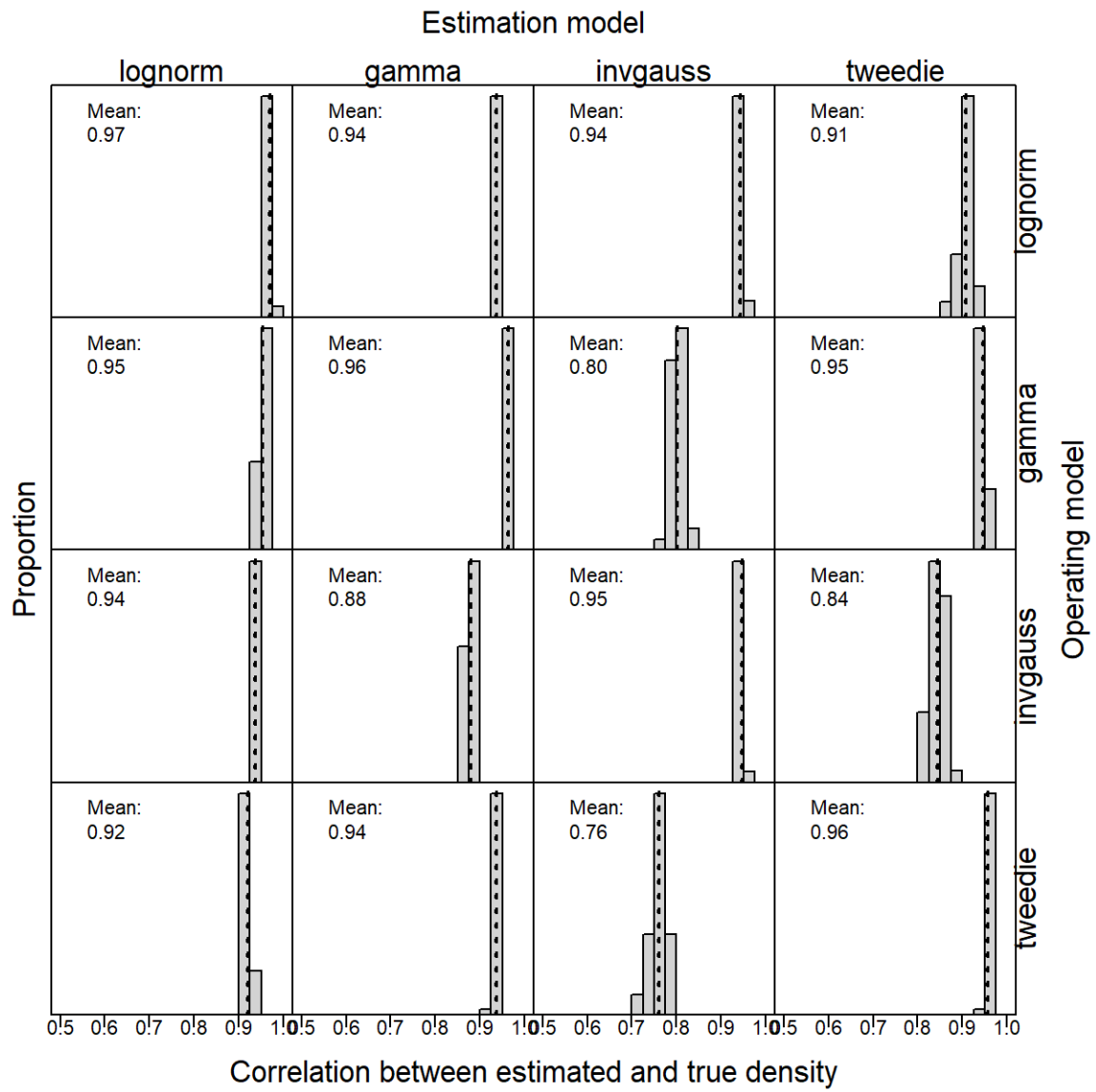
867



869

870

871 Fig. 6



872

873