# VAST model structure and user interface

## James Thorson

**Purpose of document**:

R package VAST includes many different forms of documentation, which are documented on the [package GitHub page](). This "VAST model structure and user interface" document is intended to complement these other resources by documenting and describing the model structure (all model equations and notation) while linking it to user-options that are available via the R interface to package VAST.

**Package architecture**:

VAST is developed as an R package available on GitHub. It depends upon helper functions that are bundled in package FishStatsUtils, and these helper functions are installed separately because they are also used by other spatio-temporal packages (e.g., EOFR). VAST and FishStatsUtils use S3 objects to ease interpretation of objects that are commonly saved to terminal (see Table 1 for list). VAST can be run using two primary levels of abstraction:

1. *High-level wrapper functions*: New users are recommended to explore using `FishStatsUtils::make_settings` and `FishStatsUtils::fit_model` to run VAST, and to explore results using `plot` and `summary`.

2. *Mid-level utilities*: Experienced users often run lower-level functions to accomplish basic tasks in spatial analysis, using `FishStatsUtils::make_extrapolation_info`, `FishStatsUtils::make_spatial_info`, `VAST::make_data`, and `VAST::make_model` individually.

25  Updates to VAST are released using semantic-version numbering (e.g., version 3.2.0) and a

26  battery of integrated tests (comparing results using updated code to saved results from earlier

27  versions) are run prior to numbered releases to ensure that results are backwards compatible.

28

## Model description:

30      In the following, I use mathematical notation similar to the C++ code used to define

31  the model in TMB: I use parentheses to indicate a parameter or variable that is indexed by

32  the specified indices, and I use subscripts for naming (e.g., to indicate different parameters

33  for different model components). Notation is close to common recommendations, e.g.,

34  Edwards and Auger-Méthé (2019), although I use parentheses to indicate indices of vectors,

35  matrices, and arrays, and reserve subscripts for naming. Feel free to change notation when

36  describing the model to suit your purposes in reports or publications. For further details

37  regarding terminology, motivation, and statistical properties, please read the papers listed on

38  the GitHub main page.

39

**Model Overview**

41      VAST predicts variation in density across multiple locations $s$, time intervals $t$, for

42  multiple categories $c$. Categories could include either multiple species, and/or multiple

43  size/age/sex classes for each individual species. VAST approximates the covariance between

44  these multiple categories and years using a factor-model decomposition (Thorson et al.

45  2015b, 2016a), i.e., by summing across the contribution of multiple random effects (termed

46  factors). If there is only a single category, the model reduces to a standard univariate spatio-

47  temporal model.

48      After estimating variation in density across space, time, and among categories, VAST

49  then predicts variables at extrapolation-grid cells distributed within across a user-specified

50    spatial domain.  This allows derived quantities to be calculated by summing across this

51    spatial domain; this is analogous to an "area-weighting" approach to index standardization,

52    and the resulting prediction of total abundance can be used an index of abundance.

53         In addition to spatial and spatio-temporal covariance among multiple categories,

54    VAST allows users to specify either density or catchability covariates.  Both explain variation

55    in observed catch-rate data, but VAST predicts density (for use in calculating the abundance

56    index) using density covariates but not catchability covariates.  Therefore, VAST "controls

57    for" catchability covariates when calculating an index (i.e., removes their estimated effect)

58    while "conditioning on" density covariates when calculating an index (i.e., uses them to

59    improve interpolated/extrapolated predictions of density).

60         VAST estimates the value of spatial variables at $n_x$ knots, as well as additional

61    boundary vertices such that the total number of spatial locations is $n_s$.  VAST specifically

62    uses a k-means algorithm to identify the location of $n_x$ knots to minimize the total distance

63    between the location of knots and either data or extrapolation-grid cells.  This distributes

64    knots as a function of the spatial intensity of sampling data.

65

66    **Linear predictors**

67    The model potentially includes two linear predictors (because it is designed to support delta-

68    models, which include two components).  The first linear predictor $p_1(i)$ represents

69    encounter probability in a delta-model, or zero-inflation in a count-data model:

70    $$p_1(i) = \underbrace{\mu_{\beta 1}(c_i) + \beta_1(c_i, t_i)}_{Temporal\ variation} + \underbrace{\omega_1^*(s_i, c_i)}_{Spatial\ variation} + \underbrace{\varepsilon_1^*(s_i, c_i, t_i)}_{Spatio-temporal\ variation} + \underbrace{\eta_1(v_i, c_i)}_{Vessel\ effects}$$

71    $$+ \underbrace{v_1(c_i, t_i)}_{Habitat\ covariates} + \underbrace{\zeta_1(i)}_{Catchability\ covariate}$$

72    where $p_1(i)$ is the predictor for observation $i$, arising for category $c_i$ at location $s_i$ and time

73    $t_i$.  Similarly, the second linear predictor $p_2(i)$ represents positive catch rates in a delta-

3

74   model, or the count-data intensity function in a count-data model, where all variables and

75   parameters are defined similarly except using different subscripts (Thorson and Barnett 2017;

76   Thorson 2019).  Model components are specified hierarchically to generate a correlation

77   among categories and years as explained next.

78

79   **Temporal variation**

80   Regarding intercepts representing temporal variation:

$$\beta_1(c, t) = \sum_{f=1}^{n_{\beta 1}} L_{\beta 1}(c_i, f)\beta_1(t_i, f)$$

82   where $\beta_1^*(t_i, f)$ represents temporal variation for time $t_i$ for factor $f$ (of $n_{\beta 1}$ factors

83   representing temporal variation), $L_{\beta 1}(c_i, f)$ is the loadings matrix that generates temporal

84   covariation among categories for this linear predictor, and $\beta_2^*(c_i)$ represents the time-average

85   for each category $c_i$.  The number of factors $n_{\beta 1}$ can range from zero to the number of

86   categories $n_c$, $0 \leq n_{\beta 1} \leq n_c$, where $n_{\beta 1} = 0$ is equivalent to eliminating all temporal terms

87   from the model.  By default, $n_{\beta 1} = n_c$, $\beta_1(t, f)$ is treated as a fixed effect for each year $t$ and

88   factor $f$, and $\mathbf{L}_{\beta 1}$ is an identity matrix;  this formulation is equivalent to estimating a separate

89   intercept $\beta_1(t_i, c) = \beta_1(t_i, f)$ for each category and year.  However, the intercepts can

90   instead be treated as a random effect using the factor-model formulation, which allows for

91   sharing information among years and categories.  When treated as random, $\beta_1(t_i, f)$ is

92   assigned a normal distribution with unit variance, such that $\mathbf{L}_{\beta 1}^T \mathbf{L}_{\beta 1}$ is the covariance among

93   categories for a given process (Thorson et al. 2015b).  When treating intercepts as random,

94   and when there is only one category and using one factor ($n_{\beta 1} = 1$), then $\mathbf{L}_{\beta 1}$ is a 1x1 matrix

95   (i.e. a scalar) such $\mathbf{L}_{\beta 1}^2$ is the variance and the absolute value, $abs(\mathbf{L}_{\beta 1})$ is the standard

96   deviation for temporal variation.

97    By default the model specifies that each intercept $\beta_1(c,t)$ and $\beta_2(c,t)$ is a fixed

98    effect. However, other settings specify the following autocorrelation structure:

99
$$\beta_1(t,f) \sim \begin{cases} Normal(0,1) & \text{if } t = t_{min} \\ Normal(\rho_{\beta 1}\beta_1(t-1,f),1) & \text{if } t > t_{min} \end{cases}$$

100
$$\beta_2(t,f) \sim \begin{cases} Normal(0,1) & \text{if } t = t_{min} \\ Normal(\rho_{\beta 2}\beta_2(t-1,f),1) & \text{if } t > t_{min} \end{cases}$$

101    Where $t_{min}$ is the index for the first modelled year and $\rho_{\beta 1}$ and $\rho_{\beta 2}$ are the estimated degree

102    of first-order autocorrelation in temporal variation (note that random effects have a variance

103    of one given that they are subsequently multiplied by loadings matrices that represent the

104    temporal covariance among factors). Options include:

105    1. *Independent among years* –specifies $\rho_{\beta 1} = 0$

106    2. *Random walk* –specifies $\rho_{\beta 1} = 1$

107    3. *Constant intercept* –specifies $\rho_{\beta 1} = 0$ and $\sigma_{\beta 1}^2 = 0$ (i.e., $\beta_1(t)$ is constant for all $t$)

108    4. *Autoregressive* –estimates $\rho_{\beta 1}$ as a fixed effect

109    and settings are defined identically for specifying $\rho_{\beta 2}$.

110

111    **Spatial variation**

112    Regarding spatial variation:

113
$$\omega_1^*(s,c) = \sum_{f=1}^{n_{\omega 1}} L_{\omega 1}(c_i,f)\omega_1^*(s_i,f)$$

114    where $\omega_1^*(s_i,f)$ represents predicted spatial variation in the first linear predictor occurring at

115    the location $s_i$ of sample $i$ for factor $f$ (of $n_{\omega 1}$ factors representing spatial variation), and

116    $L_{\omega 1}(c_i,f)$ is the loadings matrix that generates spatial covariation among categories for this

117    linear predictor.

118     VAST specifies internally that the spatial and spatio-temporal Gaussian random fields

119     (GMRFs) have a variance of 1.0. By default VAST estimates their values at each of $n_s$

120     vertices as follows:

121                             $$\boldsymbol{\omega}_1(f) \sim MVN(\mathbf{0}, \mathbf{R}_1)$$

122                             $$\boldsymbol{\omega}_2(f) \sim MVN(\mathbf{0}, \mathbf{R}_2)$$

123     where $\boldsymbol{\omega}_1(f)$ is the vector of length $n_s$ formed when subsetting $\omega_1(s, f)$ for a given $f$.

124     Specifying a variance of 1.0 ensures that the covariance among categories is defined by the

125     loadings matrix for that term. These GMRFs are then projected to calculate their value at

126     every location $s_i$ using matrix $\mathbf{A}$ with $n_i$ rows and $n_s$ columns. Values are then predicted as

127     e.g.:

128                             $$\boldsymbol{\omega}_1^*(f) = \mathbf{A}_i\boldsymbol{\omega}_1(f)$$

129     where $\boldsymbol{\omega}_1^*(f)$ is the vector of length $n_i$, containing the predicted value $\omega_1^*(s_i, f)$ for spatial

130     variation in the first linear predictor at every location $s_i$, and other spatial variables are

131     predicted similarly using matrix $\mathbf{A}$.

132

133     **Spatio-temporal variation**

134     Regarding spatio-temporal variation:

135                 $$\varepsilon_1^*(s, c, t) = \sum_{f=1}^{n_{\varepsilon c1}} L_{\varepsilon c1}(c_i, f)\varepsilon_1^*(s_i, f, t_i)$$

136     where $\varepsilon_1^*(s_i, f, t_i)$ represents predicted spatio-temporal variation in the first linear predictor

137     for each factor $f_1$ representing covariance among species (of $n_{\varepsilon c1}$ factors) and each factor $f_2$

138     representing covariance among years (of $n_{\varepsilon t1}$ such factors). Meanwhile and $L_{\varepsilon c1}(c_i, f_1)$ is

139     the loadings matrix that generates spatio-temporal covariation among species, and $L_{\varepsilon t1}(t_i, f_2)$

140     is the loadings matrix that generates spatio-temporal covariation among years.

141    By default, the model specifies that each vector of spatio-temporal random effects,

142    $\boldsymbol{\varepsilon}_1(f_1, f_2)$ and $\boldsymbol{\varepsilon}_2(f_1, f_2)$ composed of $\varepsilon_1(s, f_1, f_2)$ and $\varepsilon_2(s, f_1, f_2)$ across locations $s$, is

143    independent for each factor representing covariation among categories ($f_1$) and among years

144    ($f_2$). We describe the process for the 1$^{st}$ linear predictor, and an identical process is used for

145    the 2$^{nd}$ linear predictor (using different subscripts):

146    $$\boldsymbol{\varepsilon}_1(f_1, f_2) \sim MVN(\mathbf{0}, \mathbf{R}_1)$$

147    Values are then predicted as e.g.:

148    $$\boldsymbol{\varepsilon}_1^*(f_1, f_2) = \mathbf{A}_i \boldsymbol{\varepsilon}_1(f_1, f_2)$$

149    This is then projected across years and categories using loadings matrices $\mathbf{L}_{\varepsilon_t 1}$ and $\mathbf{L}_{\varepsilon_c 2}$:

150    $$\varepsilon_1'(s, c, t) = \sum_{f_1=1}^{n_{\varepsilon c1}} \sum_{f_2=1}^{n_{\varepsilon t1}} L_{\varepsilon_c 1}(c, f_1) L_{\varepsilon_t 1}(f_2, t) \varepsilon_1(s, f_1, f_2)$$

151    Using a factor-decomposition to approximate covariation among years is a generalization of

152    empirical orthogonal function (EOF) analysis (Thorson et al. 2020). The user then can

153    specify a vector-autoregressive structure:

154    $$\varepsilon_1(s, c_1, t) = \begin{cases} \varepsilon_1'(s, c_1, t) & \text{if } t = t_{min} \\ \sum_{c_2=1}^{n_c} b(c_1, c_2) \varepsilon_1'(s, c_2, t-1) & \text{if } t > t_{min} \end{cases}$$

155    Where $b(c_1, c_2)$ is the estimated impact of spatio-temporal variation in category $c_2$ on spatio-

156    temporal changes in category $c_1$:

157    $$b(c_1, c_2) = \begin{cases} \sum_{f=1}^{n_b} \chi(c_1, f) \psi(f, c_2) + \rho_{\varepsilon 1}(c_1) & \text{if } c_1 = c_2 \\ \sum_{f=1}^{n_b} \chi(c_1, f) \psi(f, c_2) & \text{if } c_1 \neq c_2 \end{cases}$$

158    Where $\chi(c_1, f)$ and $\psi(f, c_2)$ represent elements of matrices $\mathbf{X}$ and $\boldsymbol{\Psi}$, where the product $\mathbf{X}\boldsymbol{\Psi}$

159    is the typical interaction matrix in a cointegration model (Engle and Granger 1987), where $\boldsymbol{\Psi}$

160    projects dynamics to a low-dimensional subspace and $\mathbf{X}$ represents responses within that

161  subspace. By default $n_b = 0$ corresponding to $\mathbf{X\Psi} = \mathbf{0}$, and these terms drop out of the

162  model; however, they allow a parsimonious representation of species interactions (Thorson et

163  al. 2017, 2019). Meanwhile $\rho_{\varepsilon 1}(c)$ is the estimated degree of first-order autocorrelation in

164  temporal variation:

165  1. *Random walk* – specifies $\rho_{\varepsilon 1}(c) = 1$

166  2. *Autoregressive* – estimates $\rho_{\varepsilon 1}$ as a single fixed effect with the same value for all

167     categories

168  3. *Individual autoregressive* -- estimates a separate value of $\rho_{\varepsilon 1}(c)$ as a single fixed effect

169     for each category

170  and settings are defined identically for specifying $\rho_{\varepsilon 2}$.

171

172  **Overdisperison**

173  Regarding overdispersion:

$$\eta_1(v_i, c_i) = \sum_{f=1}^{n_{\eta 1}} L_1(c_i, f)\eta_1(v_i, f)$$

174

175  where $\eta_1(v_i, f)$ represents random variation in catchability among a grouping variable (tows

176  or vessels) for each factor $f$ (of $n_{\eta 1}$ factors representing overdispersion), and $L_1(c_i, f)$ is a

177  loadings matrix that generates covariation in catchability among categories for this predictor.

178  All loadings matrices are specified similarly to $\mathbf{L}_{\beta 1}$, i.e., where factors have a variance of one

179  such that $\mathbf{L}^T \mathbf{L}$ represents the covariance among categories. The main difference is that

180  spatial, spatio-temporal, and overdispersion factors can only be specified as random effects,

181  while the intercepts can be specified as either random or fixed (where specifying as fixed

182  "turns off" all factor-modelling for that intercept).

183

184  **Density covariates**

185    Regarding covariates affecting densities ("density" or "habitat" covariates):

$$v_1(c_i, t_i) = \sum_{p=1}^{n_p} \left( \gamma_1(c_i, p) + \sigma_{\xi 1}(c_i, p)\xi_1^*(s_i, c_i, p) \right) X(i, t_i, p)$$

187    where $X(i, t_i, p)$ is an three-dimensional array of $n_p$ measured density covariates that explain

188    variation in density for time $t$ and the location $s_i$ where sampling occurred for sample $i$.

189    VAST can include a separate, spatially-varying effect of each habitat covariate $p$ for each

190    category $c$. The spatially varying slope is $\gamma_1(c_i, t_i, p) + \sigma_{\xi 1}(c, p)\xi_n(s, c, p)$, where

191    $\gamma_1(c_i, t_i, p)$ is the average effect of density covariate $X(i, t_i, p)$ for category $c$, $\xi_n(s_i, c_i, p)$

192    represents spatial variation in that effect (which has a mean of zero and standard deviation of

193    one), and $\sigma_{\xi 1}(c, p)$ represents the estimated standard deviation of spatial variation of

194    covariate $p$ for category $c$. By default VAST estimates spatially-varying slope terms values

195    at each vertex as follows:

$$\boldsymbol{\xi}_1(c, p) \sim MVN(\mathbf{0}, \mathbf{R}_1)$$

197    Values are then predicted as e.g.:

$$\boldsymbol{\xi}_1^*(c, p) = \mathbf{A}_i \boldsymbol{\xi}_1(c, p)$$

199

## Catchability covariates

201    Finally, regarding covariates affecting the process of obtaining measurements ("catchability"

202    or "detectability" covariates):

$$\zeta_1(i) = \sum_{k=1}^{n_k} \left( \lambda_1(k) + \sigma_{\varphi 1}(k)\varphi_1^*(s_i, k) \right) q_1(i, k)$$

204    Where $q_1(i, k)$ is an element of matrix $\mathbf{Q}_1$ composed of $n_k$ measured catchability covariates

205    that explain variation in catchability, $\lambda_1(k)$ is the estimated impact of catchability covariates

206    for this linear predictor, $\varphi_1^*(s_i, k)$ is unit-variance spatial variation in that slope term such that

207     $\sigma_{\varphi1}(k)\varphi_1^*(s_i, k)$ has standard deviation $\sigma_{\varphi1}(k)$, where spatial variation in detectability is

208     specified as follows:

209 $$\boldsymbol{\varphi}_1(k) \sim MVN(\mathbf{0}, \mathbf{R}_1)$$

210     Values are then predicted as e.g.:

211 $$\boldsymbol{\varphi}_1^*(c, p) = \mathbf{A}_i \boldsymbol{\varphi}_1(k)$$

212

213     **Link functions and observation error distributions**

214     There are currently four options for the link function. For the latest set of options see the R

215     help documentation by typing into the R terminal `?VAST::Data_Fn`.

216     1. `ObsModel[2]=0` applies a logit-link for the first linear predictor:

217 $$r_1(i) = \text{logit}^{-1}\big(p_1(i)\big)$$

218       where $r_1(i)$ is the predictor encounter probability in a delta-model, or zero-inflation in a

219       count-data model, and $logit^{-1}(p_1(i))$ is the inverse-logit (a.k.a. logistic) function of

220       $p_1(i)$, and:

221 $$r_2(i) = a_i \times \text{log}^{-1}\big(p_2(i)\big)$$

222       where $r_2(i)$ is the predicted biomass density for positive catch rates in a delta-model or

223       mean-intensity function for a count-data model, $log^{-1}(p_2(i))$ is the exponential function

224       of $p_2(i)$, and $a_i$ is the area-swept for observation $i$, which enters as a linear offset for

225       expected biomass given an encounter.

226     2. `ObsModel[2]=1` corresponds to a "Poisson-link" delta-model that approximates a Tweedie

227       distribution:

228 $$r_1(i) = 1 - \text{exp}\big(-a_i \times \text{exp}(p_1(i))\big)$$

229       where $r_1(i)$ is the predictor encounter probability and $1 - \text{exp}\big(-a_i \times \text{exp}(p_1(i))\big)$ is a

230       complementary log-log link of $p_1(i) + \text{log}(a_i)$, and:

$$231 \qquad r_2(i) = \frac{a_i \times \exp(p_1(i))}{r_1(i)} \times \exp(p_2(i))$$

232    where $r_2(i)$ is the predicted biomass given that the species is encountered. In this

233    "Poisson-process" link function, $\exp(p_1(i))$ is interpreted as the density in number of

234    individuals per area such that $a_i \times \exp(p_1(i))$ is the predicted number of individuals

235    encountered, and $\exp(p_2(i))$ is interpreted as the average weight per individual. Area-

236    swept $a_i$ therefore enters as a linear offset for the expected number of individuals

237    encountered (Thorson 2018). This Poisson-link function should only be used for delta-

238    models, and not for count-data models, but can also be used to combine encounter, count,

239    and biomass-sampling data (see section below for details).

240

241    **Observation models**:

242    There are different user-controlled options for observation models for available sampling

243    data. I distinguish between observation models for continuous-valued data (e.g., biomass, or

244    numbers standardized to a fixed area), and observation models for count data (e.g., numbers

245    treating area-swept as an offset). However, both are parameterized such that the expectation

246    for sampling data $\mathbb{E}(B_i) = r_1(i) \times r_2(i)$.

247    *Continuous-valued data (e.g., biomass)*

248    If using an observation model with continuous support (e.g., a normal, lognormal, gamma, or

249    Tweedie models), then data $b_i$ can be any non-negative real number, $b_i \in \mathcal{R}$ and $b_i \geq 0$.

250    VAST calculates the probability of these data as:

$$251 \qquad \Pr(b_i = B) = \begin{cases} 1 - r_1(i) & \text{if } B = 0 \\ r_1(i) \times g\{B | r_2(i), \sigma_m^2(c)\} & \text{if } B > 0 \end{cases}$$

252    where `ObsModel[1]` controls the probability density function $g\{B | r_2(i), \sigma_m^2(c)\}$ used for

253    positive catch rates (see `?Data_Fn` for a list of options), where each options is defined to have

254     with expectation $r_2(i)$ and dispersion $\sigma_m^2(c)$, where dispersion parameter $\sigma_m^2(c)$ varies

255     among categories by default.

256     *Discrete-valued data (e.g., abundance)*

257     If using an observation model with discrete support (e.g., a Poisson, negative-binomial,

258     Conway-Maxwell Poisson, or lognormal-Poisson models), then data $b_i$ can be any whole

259     number, $b_i \in \{0,1,2,\dots\}$. VAST calculates the probability of these data as:

260
$$\Pr(B = b_i) = \begin{cases} \left(1 - r_1(i)\right) + g\{B = 0 | r_2(i), \dots\} & \text{if } B = 0 \\ r_1(i) \times g\{B = b_i | r_2(i), \dots\} & \text{if } B > 0 \end{cases}$$

261     where `ObsModel[1]` controls the probability mass function $g\{B | r_2(i), \dots\}$ used (again, see

262     `?Data_Fn` for a list of options), where I use … to signify that these probability mass functions

263     generally can have one or more parameter governing dispersion, and the precise number and

264     interpretation varies among observation models (i.e., the value of `ObsModel[1]`). For these

265     count-data models, $\left(1 - r_1(i)\right)$ is the "zero-inflation probability" (i.e., the proportion of

266     habitat in the immediate vicinity of location $s_i$ and time $t_i$ that is never occupied), while $r_2(i)$

267     is the expected value for probability mass function $g\{B = b_i | r_2(i), \dots\}$ (i.e., the number of

268     individuals that are in the vicinity of sampling in habitat that is occupied), and $g\{B = $

269     $0 | r_2(i), \dots\}$ is the probability of not encountering category $c$ given that sampling occurs in

270     occupied habitat (Martin et al. 2005).

271

272     **Settings regarding spatial smoothers**

273        VAST then uses a stochastic partial differential equation (SPDE) approximation to the

274     probability density function for spatial and spatio-temporal variation (Lindgren et al. 2011).

275     This SPDE approximation involves generating a triangulated mesh that has a vertex of a

276     triangle at each knot, and VAST generates this triangulated mesh using package *R-INLA*

277     (Lindgren 2012). This mesh includes all $n_x$ user-specified "interior vertices," as well as

278    additional "boundary vertices" such that the total number of interior and boundary vertices is

279    $n_s$. Outputs from this triangulated mesh can then be used to calculate the precision (inverse-

280    covariance) matrix for a multivariate normal probability density function for the value of a

281    spatial variable at all $n_s$ verticies.  Specifically, the correlation $\mathbf{R}_1(s, s + h)$ between

282    location $s$ and location $s + h$ for spatial and spatio-temporal terms included in the first linear

283    predictor is approximated as following a Matern function:

$$\mathbf{R}_1(s, s + h) = \frac{1}{2^{\nu-1}\Gamma(\nu)} \times (\kappa_1 |h\mathbf{H}|)^\nu \times K_\nu(\kappa_1 |h\mathbf{H}|)$$

285    where $\mathbf{H}$ is a two-dimensional linear transformation representing geometric anisotropy (with

286    a determinant of 1.0), $\nu$ is the Matern smoothness (fixed at 1.0), and $\kappa_1$ governs the

287    decorrelation distance for that first linear predictor ($\kappa_2$ is also separately estimated for the

288    second linear predictor).  By default, the two degrees of freedom in $\mathbf{H}$ are estimated as fixed

289    effects, but the user can specify isotropy (i.e., $\mathbf{H} = \mathbf{I}$).

290         There are also other options:

291    1.  *barrier effects*:  avoiding correlations traveling across land;

292    2.  *spherical projections*:  calculating distance based on spherical coordinates, to avoid

293        sensitivity to chosen projection;

294    3.  *stream-network distance*:  calculating distance based on river distances in a stream

295        network or other graphical spatial dependency (Hocking et al. 2018).

296

**Interpolating spatial variation from knots to the location of samples**

298    Starting with VAST release 3.0.0, users can choose between two options for smoothing

299    spatial variation.

300    1.  *Piecewise constant*:  Following the conventional for releases of VAST prior to 3.0.0,

301        users can specify `fine_scale=FALSE`. Given this specification, spatial variables at

302        location $s$ are fixed equal to their value at the nearest "knot."  This involves

303        specifying matrix $\mathbf{A}_i$ such that row $i$ has value zero except for one cell containing a

304        value of one for the knot closest to sample $i$.

305    2. *Bilinear interpolation*: Following standard practices using the software R-INLA

306        (Lindgren 2012; Lindgren and Rue 2015), users can specify `fine_scale=TRUE`. Given

307        this specification, spatial variables at location $s$ are interpolated using the triangulated

308        mesh that is also used to approximate spatial variation. Specifically, matrix $\mathbf{A}_i$ has

309        row $i$ with value zero except for three cells, representing the vertices of the triangle

310        containing location $s_i$.

311

312    **Structure on parameters among years**:

313    There are different user-controlled options for specifying structure for intercepts or spatio-

314    temporal variation across time.

315

316    **Parameter estimation**

317    Parameters are estimated using maximum likelihood, where the maximum likelihood of fixed

318    effects is obtained by integrating a joint likelihood function with respect to random effects

319    (Searle et al. 1992; Gelman and Hill 2007; Thorson and Minto 2015). This integral is

320    approximated using the Laplace approximation (Skaug and Fournier 2006), as implemented

321    in Template Model Builder (Kristensen et al. 2016). The likelihood is then optimized in the

322    R statistical environment (R Core Team 2017), and standard errors are obtained using a

323    generalization of the delta method (Kass and Steffey 1989). Derived quantities calculated via

324    a nonlinear transformation of random effects can be bias-corrected using the epsilon-method

325    (Tierney et al. 1989; Thorson and Kristensen 2016). Depending upon user-specified options,

326    different parameters will be either fixed (estimated via maximizing the log-likelihood) or

327    random (integrated across when calculating the log-likelihood). Please use R function

328 `ThorsonUtilities::list_parameters( Obj )` to see a list of estimated parameters (where `Obj` is

329 the compiled VAST object), including which are fixed or random.

330

## Combining multiple data types

332 VAST can be used to combine encounter/non-encounter, count, and biomass-sampling data.

333 This involves specifying a Poisson-link delta model which predicts each data type from

334 numbers density $\exp(p_1(i))$ and biomass-per-individual $\exp(p_2(i))$, see Grüss and Thorson

335 (2019) for details. This approach is specified by associating each observation with a given

336 error distribution using input `e_i` where e.g. `e_i[1]` is the error-distribution for the 1$^{st}$

337 observation. The user then specifies multiple observation errors via input `ObsModel_ez`:

```
338   # Control observation error
339   ObsModel_ez = cbind( "PosDist"=c(13,14,2), "Link"=c(1,1,1) )
340
```

341 In this specification, `e_i[1]==1` indicates that the first observation follows a Bernoulli

342 distribution for encounter/non-encounter data, `e_i[1]==2` indicates that this observation

343 follows a lognormal-Poisson distribution for count data, and `e_i[1]==3` indicates that it

344 follows a gamma distribution for biomass-sampling data. This specification can be modified

345 to include different combinations of these same data types.

346

## Relationship to other named models

348 VAST can be configured to be identical to (or closely mimic) many models that have

349 previously been published in ecology and fisheries:

350   1. *Spatial Gompertz model*: If intercepts are constant across years, spatio-temporal variation

351      follows an autoregressive process, and only one category is modelled, then VAST is

352      identical to a spatio-temporal Gompertz model (Thorson et al. 2014).

2. *Spatial factor analysis*: If only one year is analysed and multiple categories are modelled, VAST is similar to spatial factor analysis (Thorson et al. 2015b), although it permits the use of a delta-model (i.e., separate analysis of encounters and positive catch rates).

3. *Spatial dynamic factor analysis*: If intercepts are constant among years, spatio-temporal variation follows an autoregressive process, and multiple categories are modelled, then VAST is similar to spatial dynamic factor analysis (Thorson et al. 2016a), although VAST allows separate estimates of spatial vs. spatio-temporal covariation and also the use of a delta-model.

4. *Empirical orthogonal function analysis*: VAST can be configured to replicates empirical orthogonal function analysis, e.g., as commonly used by physical oceanographers to summarize physical conditions to produce an annual index and spatial map associated with a positive phase of the resulting index. However, I will wait to document this until the associated paper is published.

## Predicting variables across the spatial domain and calculating derived quantities

After a nonlinear minimizer has identified the value of fixed effects that maximizes the Laplace approximation to the marginal likelihood, Template Model Builder predicts the value of random effects that maximizes the joint likelihood conditional on these fixed effects. It then uses the predicted values of random effects to predict each spatial variable at each of $n_g$ "extrapolation-grid cells" that are used to summarize the spatial domain of sampling (Shelton et al. 2014; Thorson et al. 2015a). Predicting random effects at extrapolation-grid cell $g$ at location $s_g$ is accomplished using matrix $\mathbf{A}_g$ with $n_g$ rows and $n_s$ columns. Values are predicted as e.g.:

$$\boldsymbol{\omega}_1^*(f) = \mathbf{A}_g \boldsymbol{\omega}_1(f)$$

378     where $\boldsymbol{\omega}_1^*(f)$ is the vector of length $n_i$, containing the predicted value $\omega_1^*(s_g, f)$ for spatial

379     variation in the first linear predictor at every location $s_g$, and other spatial variables are

380     predicted similarly using matrix $\mathbf{A}_g$. Predicted values for random effects are then plugged

381     into the linear predictor, e.g.:

382
$$p_1(g,c,t) = \beta_1^*(c) + \underbrace{\sum_{f=1}^{n_{\beta 1}} L_{\beta 1}(c,f)\beta_1(t,f)}_{Temporal\ variation} + \underbrace{\sum_{f=1}^{n_{\omega 1}} L_{\omega 1}(x,f)\omega_1^*(g,f)}_{Spatial\ variation}$$

383
$$+ \underbrace{\sum_{f=1}^{n_{\varepsilon 1}} L_{\varepsilon 1}(c,f)\varepsilon_1^*(g,f,t)}_{Spatio-temporal\ variation} + \underbrace{\sum_{p=1}^{n_p} \Big(\gamma_1(c,t,p) + \sigma_{\xi 1}(c,p)\xi_1^*(g,c,p)\Big) X(g,t,p)}_{Habitat\ covariates}$$

384     where $p_2(g,c,t)$ is predicted similar, and these linear predictors are used in turn to predict

385     $r_1(g,c,t)$ and $r_2(g,c,t)$, where their product is predicted biomass-density $d(g,c,t)$ at every

386     extrapolation-grid cell $g$, category $c$, and time $t$.

387         By default, density is used to predict total abundance for the entire domain (or a

388     subset of the domain) for a given species:

389
$$I(c,t,l) = \sum_{x=1}^{n_x} \Big(a(g,l) \times d(g,c,t)\Big)$$

390     where $a(g,l)$ is the area associated with extrapolation-grid cell $g$ for index $l$; and. The user

391     can also specify additional post-hoc calculations via the Options vector:

```
392     Options = c("SD_site_density"=0, "SD_site_logdensity"=0, "Calculate_Range"=0,
393     "Calculate_evenness"=0, "Calculate_effective_area"=0, "Calculate_Cov_SE"=0,
394     'Calculate_Synchrony'=0, 'Calculate_Coherence'=0)
395
```

396     1. *Distribution shift* $-$ RhoConfig[3]=1 turns on calculation of the centroid of the

397         population's distribution:

398
$$Z(c,t,m) = \sum_{x=1}^{n_x} \frac{\big(z(g,m) \times a(g,1) \times d(g,c,t)\big)}{I(c,t,1)}$$

17

399    where $z(g, m)$ is a matrix representing location for each extrapolation-grid cell (by

400    default $z(g, m)$ is the location in Eastings and Northings of each knot), representing

401    movement North-South and East-West).  This model-based approach to estimating

402    distribution shift can account for differences in the spatial distribution of sampling, unlike

403    conventional sample-based estimators (Thorson et al. 2016b).

404    2. *Range expansion* – `RhoConfig[5]=1` turns on calculation of effective area occupied.  This

405    involves calculating biomass-weighted average density:

$$D(c, t, l) = \sum_{x=1}^{n_x} \frac{a(x, l) \times d(x, c, t)}{I(c, t, l)} d(x, c, t)$$

406

407    Effective area occupied is then calculated as the area required to contain the population at

408    this average density:

$$A(c, t, l) = \frac{I(c, t, l)}{D(c, t, l)}$$

409

410    This effective-area occupied estimator can then be used to monitor range expansion or

411    contraction or density-dependent range expansion (Thorson et al. 2016c).

412    The calculation of these and other derived quantities can be turned on and off using input

413    `Options` to function `make_data` (see Table 2).

414

## List of features

416    I next provide a list of "features" organized as decisions that can be made by the analyst.

417    Although this is somewhat redundant with the explanations provided above, this list might be

418    useful for some readers to provide a high-level overview of different options that are

419    available.  This "feature set" is also provided as a high-level summary of what VAST is

420    designed to be capable of doing; any software replacing VAST would ideally include this

421    same set of features.

*Basic features in a generalized linear model (GLM)*

422

423    1. Specifying one of several possible distributions for data, including for:

424        a. Count data using a Poisson, negative-binomial, Conway-Maxwell-Poisson, or

425        Poisson-lognormal distribution, including zero-inflated versions of each;

426        b. Continuous-valued data that include zeros using a delta-model with a lognormal

427        or gamma distribution for positive values.

428    2. Specifying one of several possible link functions for predicting data given linear

429    predictors including:

430        a. A conventional delta-model;

431        b. A Poisson-link delta model.

432    3. Including dynamic habitat covariates or not;

433    4. Including catchability covariates or not;

*Basic features in a spatio-temporal generalized linear mixed model (GLMM)*

434

435    5. Specify an "extrapolation grid" using input

436    `FishStatsUtils::make_extrapolation_info(..., Region)`, which is used to calculate the

437    area associated with each knot $a_x$. This can be a user-specified extrapolation grid if

438    `FishStatsUtils::make_extrapolation_info(..., Region="User", input_grid=Input)`,

439    where `Input` is a data frame supplied by the user.

440    6. Specifying a method for defining "knots";

441    7. Specifying the number of "knots";

442    8. Spatial variation being estimated ("turned on") or ignored ("turned off") for either linear

443    predictor #1 or #2;

444    9. Spatio-temporal variation being estimated ("turned on") or ignored ("turned off") for

445    either linear predictor #1 or #2;

446    10. Specifying that habitat covariates can affect linear predictors different ways including as:

447       a. a linear effect;

448       b. a spatially-varying effect; or

449       c. both linear and spatially-varying effects simultaneously.

450 *Multivariate analysis*

451 11. Including a "multivariate" structure with multiple responses that covary due to a specified

452     number of "factors" for spatial and spatio-temporal terms;

453 12. Rotate results prior to interpretation, using either:

454       a. principle components rotation; or

455       b. varimax rotation.

456 *Decisions regarding temporal structure*

457 13. Annual intercepts being structured over time, including:

458       a. estimated as fixed effects in every year;

459       b. fixed as fixed effect with the same value for all years;

460       c. estimated as a random effect with independent deviations in each year;

461       d. estimated as a random effect with first-order autoregressive structure; or

462       e. estimated as a random effect with a random-walk structure.

463 14. Spatio-temporal variation being structured over time, including:

464       a. estimated as independent deviations in each year;

465       b. estimated as following a first-order autoregressive structure over time;

466       c. estimated as following a random-walk structure over time; or

467       d. estimated as following a vector-autoregressive structure involving a matrix of 1$^{st}$

468         order autoregressive interactions.

469 *Derived quantities*

470 15. Specifying spatial strata for use when calculating derived quantities;

471 16. Calculating one of many possible "univariate derived quantities", including:

472       a. abundance indices;

473       b. range shift;

474       c. effective area occupied

475       d. covariance among categories within a multivariate model; or

476       e. synchrony among categories.

477 17. Calculating "multivariate derived quantities" that are derived from estimates for multiple

478     categories in a multivariate model, e.g., where one category represents a standardized diet

479     sample (e.g., prey biomass per predator biomass in a stomach-content sample) and

480     another category represents a biomass-density sample (e.g., predator biomass in a bottom-

481     trawl sample) such that their product represents predator-expanded consumption.

482 *Unusual circumstances and special cases*

483 18. Specifying separate distributions for different data sets (e.g., when multiple surveys

484     providing different data types are available);

485 19. Specifying that some data are predicted based on summing linear predictors across

486     multiple variables (e.g., when modelling density for different size classes, and specifying

487     that some data are aggregated measurements of multiple sizes-classes);

488 20. Specifying multiple "seasons" (e.g., when modelling data with both annual and monthly

489     spatio-temporal variation).

490

491 **Common problems**

492 There are two basic problems that are often encountered during spatio-temporal delta-

493 GLMMs:

494 1. *Encounter rates*: Some combination of categories and year has 0% or 100% encounter

495     rate. If there is 100% encounter rate for category $c$ in year $t$, then $\beta_1(c, t) \to \infty$ and/or

496     $\varepsilon_1(s, c, t) \to \infty$ for that year. If there is 0% encounter rate in year $t$, then $\beta_1(c, t) \to -\infty$

497 and/or $\varepsilon_1(s, c, t) \rightarrow -\infty$ and there is no information to estimate $\beta_2(c, t)$ or $\varepsilon_2(s, c, t)$ for

498 that category $c$ and year $t$;

499 2. *Bounds*: Some parameter(s) hits a bound;

500 These problems can be solved by:

501 1. *Encounter rates*: constraining terms that vary among years (e.g., intercept $\beta$ and spatio-

502 temporal variation $\varepsilon(s, t, p)$). This can be done in many different ways that are each

503 idiosyncratic and require some special justification. The easiest options are:

504      a. If there is a small number of years with 100% encounter rate, try `ObsModel[2]=3`.

505         This indicates that VAST should check for species-years combinations with 100%

506         encounter rates and fix corresponding intercepts for encounter probability to an

507         extremely high value.

508      b. If there is a small number of years with either 100% of 0% encounter rate, add

509         temporal structure to intercepts and spatio-temporal terms using `RhoConfig`

510         options.

511      c. Four other options are listed on the [wiki](#).

512 2. *Bounds*: Please try running the model without estimating standard errors or a final

513 newton step:

```
514     # Specify derived quantities to calculate
515     TMBhelper::fit_tmb( ..., getsd=FALSE, newtonsteps=0 )
516     Then check what parameters are being estimated near an upper or lower boundary.
```

517

## How to implement basic model changes

519 There are a few basic model types that users often want to fit using VAST. I briefly describe

520 how these can be done here.

1. *Fitting encounter/non-encounter data*:  If the user wishes to use only the first component of a delta-model, i.e., to fit a binomial model to simply predict encounter probabilities, then, the `ObsModel` vector should be set to `c("PosDist"=[Make Choice], "Link"=0)`, where [Make Choice] can be any option for continuous data (i.e., 0, 1, or 2).  The user should then turn off the last two elements of the `FieldConfig` vector (i.e., `FieldConfig[3]=0` and `FieldConfig[4]=0`) such that there is no spatial or spatio-temporal variability in positive catch rates, and also turn off annual variation in the intercept for positive catch rates (i.e., `RhoConfig[2]=3`).  Finally, the user should "jitter" their presence observations by a very small amount (i.e., add a random normal deviation with a very small standard deviation, `rnorm(n=1,mean=0,sd=0.001)`, to each observation for which `b_i=1`).  This will result in VAST estimating a logistic regression model for encounter/non-encounter data, except with one additional parameter estimated ($\sigma_M$), plus one additional parameter per category ($\beta_2(c)$), where these additional parameters have no impact on other parameters, are not meant to be interpreted statistically or biologically, and are an artefact of using VAST (which is designed to fit a delta-model) to encounter/non-encounter data.  This feature has been used to estimate species distributions for use in ecosystem models (Grüss et al. 2017, 2018).

## Acknowledgements

546　NOAA scientists who have served on sampling vessels that provided data to test these

547　methods.　Finally, I think A. Grüss and S. Hoyle for providing edits to this document.

548

549

# Works cited

Edwards, A.M., and Auger-Méthé, M. 2019. Some guidance on using mathematical notation in ecology. Methods Ecol. Evol. **10**(1): 92–99. doi:10.1111/2041-210X.13105.

Engle, R.F., and Granger, C.W. 1987. Co-integration and error correction: representation, estimation, and testing. Econom. J. Econom. Soc.: 251–276.

Gelman, A., and Hill, J. 2007. Data analysis using regression and multilevel/hierarchical models. Cambridge University Press, Cambridge, UK.

Godefroid, M., Boldt, J.L., Thorson, J.T., Forrest, R., Gauthier, S., Flostrand, L., Ian Perry, R., Ross, A.R.S., and Galbraith, M. 2019. Spatio-temporal models provide new insights on the biotic and abiotic drivers shaping Pacific Herring (Clupea pallasi) distribution. Prog. Oceanogr. **178**: 102198. doi:10.1016/j.pocean.2019.102198.

Grüss, A., and Thorson, J.T. 2019. Developing spatio-temporal models using multiple data types for evaluating population trends and habitat usage. ICES J. Mar. Sci. **76**(6): 1748–1761. doi:10.1093/icesjms/fsz075.

Grüss, A., Thorson, J.T., Babcock, E.A., and Tarnecki, J.H. 2018. Producing distribution maps for informing ecosystem-based fisheries management using a comprehensive survey database and spatio-temporal models. ICES J. Mar. Sci. **75**(1): 158–177. doi:10.1093/icesjms/fsx120.

Grüss, A., Thorson, J.T., Sagarese, S.R., Babcock, E.A., Karnauskas, M., Walter, J.F., and Drexler, M. 2017. Ontogenetic spatial distributions of red grouper (Epinephelus morio) and gag grouper (Mycteroperca microlepis) in the U.S. Gulf of Mexico. Fish. Res. **193**(Supplement C): 129–142. doi:10.1016/j.fishres.2017.04.006.

Hocking, D.J., Thorson, J.T., O'Neil, K., and Letcher, B.H. 2018. A geostatistical state-space model of animal densities for stream networks. Ecol. Appl. **28**(7): 1782–1796. doi:10.1002/eap.1767.

Kass, R.E., and Steffey, D. 1989. Approximate Bayesian inference in conditionally independent hierarchical models (parametric empirical bayes models). J. Am. Stat. Assoc. **84**(407): 717–726. doi:10.2307/2289653.

Kristensen, K., Nielsen, A., Berg, C.W., Skaug, H., and Bell, B.M. 2016. TMB: Automatic differentiation and Laplace approximation. J. Stat. Softw. **70**(5): 1–21. doi:10.18637/jss.v070.i05.

Lindgren. 2012. Continuous domain spatial models in R-INLA. ISBA Bull. **19**(4): 14–20.

Lindgren, F., and Rue, H. 2015. Bayesian spatial modelling with r-inla. J. Stat. Softw. **63**(19): 1–25. doi:10.18637/jss.v063.i19.

Lindgren, Rue, H., and Lindström, J. 2011. An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach. J. R. Stat. Soc. Ser. B Stat. Methodol. **73**(4): 423–498. doi:10.1111/j.1467-9868.2011.00777.x.

Martin, T.G., Wintle, B.A., Rhodes, J.R., Kuhnert, P.M., Field, S.A., Low-Choy, S.J., Tyre, A.J., and Possingham, H.P. 2005. Zero tolerance ecology: improving ecological inference by modelling the source of zero observations. Ecol. Lett. **8**(11): 1235–1246.

R Core Team. 2017. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria. Available from https://www.R-project.org/.

Searle, S.R., Casella, G., and McCulloch, C.E. 1992. Variance components. John Wiley & Sons, Hoboken, New Jersey.

Shelton, A.O., Thorson, J.T., Ward, E.J., and Feist, B.E. 2014. Spatial semiparametric models improve estimates of species abundance and distribution. Can. J. Fish. Aquat. Sci. **71**(11): 1655–1666. doi:10.1139/cjfas-2013-0508.

599 Skaug, H., and Fournier, D. 2006. Automatic approximation of the marginal likelihood in
600     non-Gaussian hierarchical models. Comput. Stat. Data Anal. **51**(2): 699–709.
601 Thorson, J.T. 2018. Three problems with the conventional delta-model for biomass sampling
602     data, and a computationally efficient alternative. Can. J. Fish. Aquat. Sci. **75**(9):
603     1369–1382. doi:10.1139/cjfas-2017-0266.
604 Thorson, J.T. 2019. Guidance for decisions using the Vector Autoregressive Spatio-Temporal
605     (VAST) package in stock, ecosystem, habitat and climate assessments. Fish. Res. **210**:
606     143–161. doi:10.1016/j.fishres.2018.10.013.
607 Thorson, J.T., Adams, G., and Holsman, K. 2019. Spatio-temporal models of intermediate
608     complexity for ecosystem assessments: A new tool for spatial fisheries management.
609     Fish Fish. **20**(6): 1083–1099. doi:10.1111/faf.12398.
610 Thorson, J.T., and Barnett, L.A.K. 2017. Comparing estimates of abundance trends and
611     distribution shifts using single- and multispecies models of fishes and biogenic
612     habitat. ICES J. Mar. Sci. **74**(5): 1311–1321. doi:10.1093/icesjms/fsw193.
613 Thorson, J.T., Ciannelli, L., and Litzow, M.A. 2020. Defining indices of ecosystem
614     variability using biological samples of fish communities: A generalization of
615     empirical orthogonal functions. Prog. Oceanogr. **181**: 102244.
616     doi:10.1016/j.pocean.2019.102244.
617 Thorson, J.T., and Haltuch, M.A. 2018. Spatiotemporal analysis of compositional data:
618     increased precision and improved workflow using model-based inputs to stock
619     assessment. Can. J. Fish. Aquat. Sci. **76**(3): 401–414. doi:10.1139/cjfas-2018-0015.
620 Thorson, J.T., Ianelli, J.N., Larsen, E.A., Ries, L., Scheuerell, M.D., Szuwalski, C., and
621     Zipkin, E.F. 2016a. Joint dynamic species distribution models: a tool for community
622     ordination and spatio-temporal monitoring. Glob. Ecol. Biogeogr. **25**(9): 1144–1158.
623     doi:10.1111/geb.12464.
624 Thorson, J.T., and Kristensen, K. 2016. Implementing a generic method for bias correction in
625     statistical models using random effects, with spatial and population dynamics
626     examples. Fish. Res. **175**: 66–74. doi:10.1016/j.fishres.2015.11.016.
627 Thorson, J.T., and Minto, C. 2015. Mixed effects: a unifying framework for statistical
628     modelling in fisheries biology. ICES J. Mar. Sci. J. Cons. **72**(5): 1245–1256.
629     doi:10.1093/icesjms/fsu213.
630 Thorson, J.T., Munch, S.B., and Swain, D.P. 2017. Estimating partial regulation in
631     spatiotemporal models of community dynamics. Ecology **98**(5): 1277–1289.
632     doi:10.1002/ecy.1760.
633 Thorson, J.T., Pinsky, M.L., and Ward, E.J. 2016b. Model-based inference for estimating
634     shifts in species distribution, area occupied and centre of gravity. Methods Ecol. Evol.
635     **7**(8): 990–1002. doi:10.1111/2041-210X.12567.
636 Thorson, J.T., Rindorf, A., Gao, J., Hanselman, D.H., and Winker, H. 2016c. Density-
637     dependent changes in effective area occupied for sea-bottom-associated marine fishes.
638     Proc R Soc B **283**(1840): 20161853. doi:10.1098/rspb.2016.1853.
639 Thorson, J.T., Scheuerell, M.D., Olden, J.D., and Schindler, D.E. 2018. Spatial heterogeneity
640     contributes more to portfolio effects than species variability in bottom-associated
641     marine fishes. Proc R Soc B **285**(1888): 20180915. doi:10.1098/rspb.2018.0915.
642 Thorson, J.T., Shelton, A.O., Ward, E.J., and Skaug, H.J. 2015a. Geostatistical delta-
643     generalized linear mixed models improve precision for estimated abundance indices
644     for West Coast groundfishes. ICES J. Mar. Sci. J. Cons. **72**(5): 1297–1310.
645     doi:10.1093/icesjms/fsu243.
646 Thorson, J.T., Skaug, H.J., Kristensen, K., Shelton, A.O., Ward, E.J., Harms, J.H., and
647     Benante, J.A. 2014. The importance of spatial models for estimating the strength of
648     density dependence. Ecology **96**(5): 1202–1212. doi:10.1890/14-0739.1.

649    Thorson, Scheuerell, M.D., Shelton, A.O., See, K.E., Skaug, H.J., and Kristensen, K. 2015b.
650        Spatial factor analysis: a new tool for estimating joint species distributions and
651        correlations in species range. Methods Ecol. Evol. **6**(6): 627–637. doi:10.1111/2041-
652        210X.12359.
653    Tierney, L., Kass, R.E., and Kadane, J.B. 1989. Fully exponential Laplace approximations to
654        expectations and variances of nonpositive functions. J. Am. Stat. Assoc. **84**(407):
655        710–716.
656
657

658 Table 1 – List of S3 objects defined in package VAST (or its primary dependency FishStatsUtils), listing S3 methods defined for each class as

659 well as the intended purpose of each method.

| S3 object | S3 methods | Purpose |
|---|---|---|
| `VAST::make_data` | `print` | De-clutter terminal output |
| `VAST::make_model` | `print` | De-clutter terminal output |
| `FishStatsUtils::make_extrapolation_info` | `print` | De-clutter terminal output |
| | `plot` | Simple organization for plotting options |
| `FishStatsUtils::make_spatial_info` | `print` | De-clutter terminal output |
| | `print` | Simple organization for plotting options |
| `FishStatsUtils::fit_model` | `print` | De-clutter terminal output |
| | `plot` | Simple organization for plotting options |
| | `summary` | Interface to access derived quantities that users may want |

660

661 Table 2 – List of slots in vector `Options`.

| Slot name | What it does | Examples of usage |
|---|---|---|
| SD_site_density | Calculate standard error for each knot or extrapolation-grid cell (very slow with fine_scale=TRUE) | - |
| SD_site_logdensity | Calculate standard error for each knot or extrapolation-grid cell (very slow with fine_scale=TRUE) | - |
| Calculate_Range | Calculate center of gravity for use in measuring distribution shifts | (Thorson et al. 2016b) |
| SD_observation_density | Calculate standard error for response variable measured in each sample | - |
| Calculate_effective_area | Calculate effective area occupied for use in measuring range expansion/contraction | (Thorson et al. 2016c) |
| Calculate_Cov_SE | Calculate standard error for correlation / covariance among categories | (Godefroid et al. 2019) |
| Calculate_Synchrony | Calculate reduction in variance associated with asynchrony among species and/or locations | (Thorson et al. 2018) |
| Calculate_Coherence | Calculate the Gini coefficient for axes of covariation among categories | - |
| Calculate_proportion | Convert indices to a proportion in a multivariate model;  breaks | (Thorson and Haltuch 2018) |

| | | |
|---|---|---|
| | separability across categories and therefore users typically use an approximation to calculate input sample size | |
| normalize_GMRF_in_CPP | Option to potentially speed up GMRF calculations, although early testing didn't indicate substantial improvements (could be explored more) | - |
| Calculate_Fratio | Option to calculate exploitation rate using a MICE-in-space model | (Thorson et al. 2019) |
| Estimate_B0 | Option to calculate unfished density using a MICE-in-space model | (Thorson et al. 2019) |
| Project_factors | Project factors to extrapolation-grid; Useful to visualize factors when using fine_scale=TRUE | - |
| treat_nonencounter_as_zero | Option to internally track combinations of category and year that are never encountered and therefore should be treated as having zero abundance; Useful for compositional-expansion | (Thorson and Haltuch 2018) |
| simulate_random_effects | Option governing behaviour of bootstrap simulator; determines whether simulator re-simulates | (Thorson et al. 2019) |

| | random-effects conditional on fixed effects (default) or not. |
|---|---|

662

663