1 # VAST user manual

2 ## James Thorson

3

4 **Purpose of document**:

5 This document is intended to document the model structure and user-options available in

6 package VAST. For guidance and examples of how to use the model, please see the

7 Rmarkdown tutorials in the GitHub "/examples" directory. In the following, I try to use

8 notation similar to the TMB code: I use parentheses to indicate a parameter or variable that is

9 indexed by the specified indices, and I use subscripts for naming (e.g., to indicate different

10 parameters for different model components). Feel free to change notation when describing

11 the model to suite your purposes. For further details regarding terminology, motivation, and

12 statistical properties, please read the papers listed on the GitHub main page.

13 **Model description**:

14 **Linear predictors**

15 I use a delta-model that includes two linear predictors. The linear predictor for encounter

16 probability:

17
$$p_1(i) = \beta_1(c_i, t_i) + \sum_{f=1}^{n_{\omega 1}} L_{\omega 1}(c_i, f)\omega_1(s_i, f) + \sum_{f=1}^{n_{\varepsilon 1}} L_{\varepsilon 1}(c_i, f)\varepsilon_1(s_i, f, t_i)$$

18
$$+ \sum_{f=1}^{n_{\delta 1}} L_{\delta 1}(v_i, f)\delta_1(v_i, f) + \sum_{p=1}^{n_p} \gamma_1(c_i, t_i, p)X(x_i, t_i, p) + \sum_{k=1}^{n_k} \lambda_1(k) Q(i, k)$$

19 where $p_1(i)$ is the predictor for observation $i$, $\beta_1(c_i, t_i)$ is an intercept for category $c_i$ and

20 year $t_i$, $\omega_1(s_i, f)$ represents spatial variation at location $s_i$ for factor $f$ and $L_{\omega 1}(c_i, f)$ is the

21 loadings matrix that generates spatial covariation among categories for this linear predictor,

22 $\varepsilon_1(s_i, f, t_i)$ is spatio-temporal variation and $L_{\varepsilon 1}(c_i, f)$ is the loadings matrix that generates

23     spatio-temporal covariation for this predictor, $\delta_1(v_i, f)$ is random variation in catchability

24     among a grouping variable (tows or vessels) and $L_{\delta 1}(v_i, f)$ is a loadings matrix that generates

25     covariation in catchability among categories for this predictor, $X(x_i, t_i, p)$ are measured

26     density covariates that explain variation in density and $\gamma_1(c_i, t_i, p)$ is the estimated impact of

27     density covariates, and $Q(i, k)$ are measured catchability covariates that explain variation in

28     catchability and $\lambda_1(k)$ is the estimated impact of catchability covariates for this linear

29     predictor. Similarly, the linear predictor for positive catch rates:

30
$$p_2(i) = \beta_2(c_i, t_i) + \sum_{f=1}^{n_{\omega 1}} L_{\omega 2}(c_i, f)\omega_1(s_i, f) + \sum_{f=1}^{n_{\varepsilon 1}} L_{\varepsilon 2}(c_i, f)\varepsilon_2(s_i, f, t_i)$$

31
$$+ \sum_{f=1}^{n_{\delta 1}} L_{\delta 2}(v_i, f)\delta_2(v_i, f) + \sum_{p=1}^{n_p} \gamma_2(c_i, t_i, p)X(x_i, t_i, p) + \sum_{k=1}^{n_k} \lambda_2(k) Q(i, k)$$

32     where all variables and parameters are defined similarly except using different subscripts

33     (Thorson and Barnett In press, Thorson et al. In press). The loadings matrices are designed

34     such that $\mathbf{L}^T\mathbf{L}$ is the covariance among categories for a given spatial or spatio-temporal

35     process (Thorson et al. 2015a), and when there is only one category $\mathbf{L}$ is a 1x1 matrix (i.e. a

36     scalar) such that its absolute value is the standard deviation for a given process. This model

37     therefore reduces to a single-species spatio-temporal model (e.g., Thorson et al. 2015b) when

38     only one category is available.

39         The user controls the number of spatial and spatio-temporal factors used for each

40     component via input:

```
41    # Control number of factors
42    FieldConfig = c("Omega1"=1, "Epsilon1"=1, "Omega2"=1, "Epsilon2"=1)
43
```

44     where `FieldConfig[1]` controls $n_{\omega 1}$, `FieldConfig[2]` controls $n_{\varepsilon 1}$, `FieldConfig[3]` controls

45     $n_{\omega 2}$, and `FieldConfig[4]` controls $n_{\varepsilon 2}$, and a value of zero "turns off" that component of

46    spatial or spatio-temporal covariation.  The user controls the number of catchability factors

47    used for each component via input:

```
48    # Control number of spatial and spatio-temporal factors
49    OverdispersionConfig = c("Delta1"=0, "Delta2"=0)
50
```

51    where `OverdispersionConfig[1]` controls $n_{\delta 1}$, and `OverdispersionConfig[2]` controls $n_{\delta 2}$,

52    and a value of zero again "turns off" that component of random covariation in catchability.

53    For example, if the user inputs:

```
54    # Control number of spatial and spatio-temporal factors
55    OverdispersionConfig = c("Delta1"=1, "Delta2"=1)
56
```

57    then there will be one random effect estimated for each unique level of `Data_Geostat$Vessel`

58    for both the first and second linear predictors.

59    **Link functions**

60    There are different user-controlled options for link-functions that calculate expected

61    encounter probability and positive catch rates given these two linear predictors.

```
62    # Control observation error
63    ObsModel = c("PosDist"=2, "Link"=0)
64
```

65    where the 2$^{\text{nd}}$ element of this vector controls the link functions.

66    1.  `ObsModel[2]=0` corresponds to a conventional delta-model:

$$r_1(i) = logit^{-1}(p_1(i))$$

68    where $r_1(i)$ is the predictor encounter probability and $logit^{-1}(p_1(i))$ is the logistic

69    function of $p_1(i)$, and:

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

71    where $r_2(i)$ is the predicted biomass density for positive catch rates, $log^{-1}(p_2(i))$ is the

72    exponential function of $p_2(i)$, and $a_i$ is the area-swept for observation $i$, which enters as a

73    linear offset for expected biomass given an encounter.

74    2.  Alternatively, `ObsModel[2]=1` corresponds to a "Poisson-process" link function that

75        approximates a Tweedie distribution:

76
$$r_1(i) = 1 - \exp\big(-a_i \times \exp(p_1(i))\big)$$

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

78        complementary log-log link of $p_1(i) + \log(a_i)$, and:

79
$$r_2(i) = \frac{a_i \times \exp(p_1(i))}{r_1(i)} \times \exp(p_2(i))$$

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

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

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

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

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

85        encountered (Thorson In review).

86    **Observation models**:

87    There are different user-controlled options for observation models for positive catch rates.

```
88    # Control observation error
89    ObsModel = c("PosDist"=2, "Link"=0)
90
```

91    VAST then calculates the probability of data as:

92
$$\Pr(b_i = B) = \begin{cases} 1 - r_1(i) & \text{if } B = 0 \\ r_1(x_i, c_i, t_i) \times g\{B | r_2(i), \sigma_m^2(c)\} & \text{if } B > 0 \end{cases}$$

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

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

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

96    among categories by default.

97    **Settings regarding spatial domain**

98    VAST approximates spatial and spatio-temporal variation as being piecewise-constant.  To

99    do so, the user specifies n_x:

```
100   # Number of knots
101   n_x = 1000
102
```

103   VAST then uses a k-means algorithm to identify the location of n_x knots to minimize the

104   total distance between the location of available data and the location of the nearest knot.  This

105   distributes knots as a function of the spatial intensity of sampling data.

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

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

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

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

110   (Lindgren 2012).  Outputs from this triangulated mesh can then be used to calculate the

111   precision (inverse-covariance) matrix for a multivariate normal probability density function

112   for the value of a spatial variable at each mesh vertex.  Specifically, the correlation

113   $\mathbf{R}_1(s, s + h)$ between location $s$ and location $s + h$ for spatial and spatio-temporal terms

114   included in the first linear predictor is approximated as following a Matern function:

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

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

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

118   distance for that first linear predictor ($\kappa_2$ is also separately estimated for the second linear predictor).

119   By default, the two degrees of freedom in $\mathbf{H}$ are estimated as fixed effects, but the user can specify

120   isotropy (i.e., $\mathbf{H} = \mathbf{I}$) by specifying:

```
121   # Turn of geometric anisotropy
122   Data = Data_Fn( …, Aniso=FALSE )
123
```

124    VAST then specifies that the spatial and spatio-temporal Gaussian random fields each

125    have a variance of 1.0. By default VAST specifies these as follows:

126    $$\omega_1(\cdot, f) \sim MVN(\mathbf{0}, \sigma_{\omega1}^2 \mathbf{R}_1)$$

127    $$\omega_2(\cdot, f) \sim MVN(\mathbf{0}, \sigma_{\omega1}^2 \mathbf{R}_2)$$

128    $$\varepsilon_1(\cdot, f, t) \sim MVN(\mathbf{0}, \sigma_{\varepsilon1}^2 \mathbf{R}_1)$$

129    $$\varepsilon_2(\cdot, f, t) \sim MVN(\mathbf{0}, \sigma_{\varepsilon2}^2 \mathbf{R}_2)$$

130    where $\omega_1(\cdot, f)$ is the vector formed when subsetting $\omega_1(s, f)$ for a given $f$, and $\sigma_{\omega1}^2$ is the

131    variance of $\omega_1(s, f)$, where other parameters are defined similarly. Specifying a variance of

132    1.0 ensures that the covariance among categories is defined by the loadings matrix for that

133    term. However, VAST allows spatio-temporal variance to be specified differently as

134    discussed in the section titled "Structure on parameters among years".

135    **Structure on parameters among years**:

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

137    temporal variation across time, using input:

```
138    # Control autoregressive structure for parameters over time
139    RhoConfig = c("Beta1"=0, "Beta2"=0, "Epsilon1"=0, "Epsilon2"=0)
140
```

141    By default (when `RhoConfig[1]=0` and `RhoConfig[2]=0`) the model specifies that each

142    intercept $\beta_1(t)$ and $\beta_2(t)$ is a fixed effect. However, other settings specify the following

143    structure:

144    $$\beta_1(t + 1) \sim Normal(\rho_{\beta1}\beta_1(t), \sigma_{\beta1}^2)$$

145    $$\beta_2(t + 1) \sim Normal(\rho_{\beta2}\beta_2(t), \sigma_{\beta2}^2)$$

146    where `RhoConfig[1]` controls the specification of $\rho_{\beta1}$:

147    1. *Independent among years* – `RhoConfig[1]=1` specifies $\rho_{\beta1} = 0$

148    2. *Random walk* – `RhoConfig[1]=2` specifies $\rho_{\beta1} = 1$

149  3. *Constant intercept* – `RhoConfig[1]=3` specifies $\rho_{\beta 1} = 0$ and $\sigma_{\beta 1}^2 = 0$ (i.e., $\beta_1(t)$ is

150     constant for all $t$)

151  4. *Autoregressive* – `RhoConfig[1]=4` estimates $\rho_{\beta 1}$ as a fixed effect

152  and settings are defined identically for `RhoConfig[2]` specifying $\rho_{\beta 2}$.

153  By default (when `RhoConfig[3]=0` and `RhoConfig[4]=0`) the model specifies that each spatio-

154  temporal random effect $\varepsilon_1(s, f, t)$ and $\varepsilon_2(s, f, t)$ is independent among years. However,

155  other settings specify the following structure

$$\varepsilon_1(s, f, t + 1) \sim MVN(\rho_{\varepsilon 1}\varepsilon_1(s, f, t), \sigma_{\varepsilon 1}^2 \mathbf{R}_1)$$

$$\varepsilon_2(s, f, t + 1) \sim MVN(\rho_{\varepsilon 1}\varepsilon_2(s, f, t), \sigma_{\varepsilon 2}^2 \mathbf{R}_2)$$

158  where `RhoConfig[3]` controls the specification of $\rho_{\varepsilon 1}$:

159  1. *Random walk* – `RhoConfig[3]=2` specifies $\rho_{\varepsilon 1} = 1$

160  2. *Autoregressive* – `RhoConfig[3]=4` estimates $\rho_{\varepsilon 1}$ as a fixed effect

161  and settings are defined identically for `RhoConfig[4]` specifying $\rho_{\varepsilon 2}$.

162  **Relationship to other named models**

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

164  previously been published in ecology and fisheries:

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

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

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

168  2. *Spatial factor analysis*: If only one year is analysed and multiple category are modelled,

169     VAST is similar to spatial factor analysis (Thorson et al. 2015a), although it permits the

170     use of a delta-model (separate analysis of encounters and positive catch rates).

171  3. *Spatial dynamic factor analysis*: If intercepts are constant among years, spatio-temporal

172     variation follows an autoregressive process, and multiple category are modelled, then

173     VAST is similar to spatial dynamic factor analysis (Thorson et al. 2016a), although

174       VAST allows separate estimates of spatial vs. spatio-temporal covariation and also the

175       user of a delta-model.

176    **Settings regarding derived quantities**

177    After a nonlinear minimizer has identified the value of fixed effects that maximizes the

178    Laplace approximation to the marginal likelihood, Template Model Builder predicts the value

179    of random effects that maximizes the joint likelihood conditional on these fixed effects.

180    Estimated values of fixed and random effects are then used to predict density $d(x, c, t)$ for :

181
$$d(x, c, t) = r_1^*(x, c, t) \times r_2^*(x, c, t)$$

182    where $r_1^*(x, c, t)$ and $r_2^*(x, c, t)$ are identical to the values specified previously, except that

183    catchability variables are excluded from their computation (i.e., $\delta_1(v, f) = 0$ and $\lambda_1(k) = 0$,

184    etc.)

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

186    subset of the domain) for a given species:

187
$$I(c, t, l) = \sum_{x=1}^{n_x} \big( a(x, l) \times d(x, c, t) \big)$$

188    where $a(x, l)$ is the area associated with extrapolation-cell $x$ for index $l$ (Shelton et al. 2014,

189    Thorson et al. 2015b).  The user can also specify additional post-hoc calculations via input:

```
190   # Control observation error
191   RhoConfig = c("SD_site_density"=0, "SD_site_logdensity"=0, "Calculate_Range"=0,
192   "Calculate_evenness"=0, "Calculate_effective_area"=0, "Calculate_Cov_SE"=0,
193   'Calculate_Synchrony'=0, 'Calculate_Coherence'=0)
194
```

195    1.  *Distribution shift* − RhoConfig[3]=1 turns on calculation of the centroid of the

196       population's distribution:

197
$$Z(c, t, m) = \sum_{x=1}^{n_x} \frac{\big( z(x, m) \times a(x, 1) \times d(x, c, t) \big)}{I(c, t, 1)}$$

198    where $z(x, m)$ is a matrix representing location for each knot (by default $z(x, m)$ is the

199    location in Eastings and Northings of each knot), representing movement North-South

200    and East-West).  This model-based approach to estimating distribution shift can account

201    for differences in the spatial distribution of sampling, unlike conventional sample-based

202    estimators (Thorson et al. 2016b).

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

204    involves calculating biomass-weighted average density:

205
$$D(c, t, l) = \sum_{x=1}^{n_x} \frac{a(x, l) \times d(x, c, t)}{I(c, t, l)} d(x, c, t)$$

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

207    this average density:

208
$$A(c, t, l) = \frac{I(c, t, l)}{\bar{d}(c, t, l)}$$

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

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

## Works cited

211

212    Lindgren, F. 2012. Continuous domain spatial models in R-INLA. ISBA Bull. **19**(4): 14–20.
213    Lindgren, F., Rue, H., and Lindström, J. 2011. An explicit link between Gaussian fields and
214        Gaussian Markov random fields: the stochastic partial differential equation approach.
215        J. R. Stat. Soc. Ser. B Stat. Methodol. **73**(4): 423–498. doi:10.1111/j.1467-
216        9868.2011.00777.x.
217    Shelton, A.O., Thorson, J.T., Ward, E.J., and Feist, B.E. 2014. Spatial semiparametric models
218        improve estimates of species abundance and distribution. Can. J. Fish. Aquat. Sci.
219        **71**(11): 1655–1666. doi:10.1139/cjfas-2013-0508.
220    Thorson, J.T. In review. Three problems with the conventional delta-model for biomass
221        sampling data, and a computationally efficient alternative.
222    Thorson, J.T., and Barnett, L.A.K. In press. Comparing estimates of abundance trends and
223        distribution shifts using single- and multispecies models of fishes and biogenic
224        habitat. ICES J. Mar. Sci. doi:10.1093/icesjms/fsw193.
225    Thorson, J.T., Ianelli, J.N., and Kotwicki, S. In press. The relative influence of temperature
226        and size structure on fish distribution shifts: a case study on walleye pollock in the
227        Bering Sea. Fish Fish.
228    Thorson, J.T., Ianelli, J.N., Larsen, E.A., Ries, L., Scheuerell, M.D., Szuwalski, C., and
229        Zipkin, E.F. 2016a. Joint dynamic species distribution models: a tool for community

230    ordination and spatio-temporal monitoring. Glob. Ecol. Biogeogr. **25**(9): 1144–1158.
231        doi:10.1111/geb.12464.
232 Thorson, J.T., Pinsky, M.L., and Ward, E.J. 2016b. Model-based inference for estimating
233        shifts in species distribution, area occupied and centre of gravity. Methods Ecol. Evol.
234        **7**(8): 990–1002. doi:10.1111/2041-210X.12567.
235 Thorson, J.T., Rindorf, A., Gao, J., Hanselman, D.H., and Winker, H. 2016c. Density-
236        dependent changes in effective area occupied for sea-bottom-associated marine fishes.
237        Proc R Soc B **283**(1840): 20161853. doi:10.1098/rspb.2016.1853.
238 Thorson, J.T., Scheuerell, M.D., Shelton, A.O., See, K.E., Skaug, H.J., and Kristensen, K.
239        2015a. Spatial factor analysis: a new tool for estimating joint species distributions and
240        correlations in species range. Methods Ecol. Evol. **6**(6): 627–637. doi:10.1111/2041-
241        210X.12359.
242 Thorson, J.T., Shelton, A.O., Ward, E.J., and Skaug, H.J. 2015b. Geostatistical delta-
243        generalized linear mixed models improve precision for estimated abundance indices
244        for West Coast groundfishes. ICES J. Mar. Sci. J. Cons. **72**(5): 1297–1310.
245        doi:10.1093/icesjms/fsu243.
246 Thorson, J.T., Skaug, H.J., Kristensen, K., Shelton, A.O., Ward, E.J., Harms, J.H., and
247        Benante, J.A. 2014. The importance of spatial models for estimating the strength of
248        density dependence. Ecology **96**(5): 1202–1212. doi:10.1890/14-0739.1.
249
250