Edinburgh Research Explorer

# Spatiotemporal Models of an Estuarine Fish Species to Identify Patterns and Factors Impacting Their Distribution and Abundance

# Estuaries and Coasts

## Spatiotemporal models of an estuarine fish species to identify patterns and factors impacting their distribution and abundance

--Manuscript Draft--

1        Spatiotemporal models of an estuarine fish species to identify patterns and factors impacting

2        their distribution and abundance

3

4        Leo Polansky[1,*], Ken B. Newman[2], Matthew L. Nobriga[1], Lara Mitchell[2]

5

6   [1]U.S. Fish and Wildlife Service, Bay Delta Fish and Wildlife Office, Sacramento, CA 95814,

7   USA

8

9   [2]U.S. Fish and Wildlife Service, Lodi Fish and Wildlife Office, Lodi CA 95240, USA

10

11   [*]Corresponding author email: leo_polansky@fws.gov

12

13   Running page head: Distribution and abundance of delta smelt

14    ABSTRACT: Understanding the distribution and abundance of organisms can be exceedingly

15    difficult for pelagic fish species that live in estuarine environments. This is particularly so for

16    fish that cannot be readily marked and released or otherwise tracked, such as the diminutive delta

17    smelt, *Hypomesus transpacificus*, endemic to the San Francisco Estuary. The environmental

18    factors that influence distribution operate at multiple scales, from daily tidal cycles and local

19    perceptual fields to seasonal and annual changes in dominant environmental gradients spanning

20    the entire San Francisco Estuary. To quantify scale specific patterns and factors shaping the

21    spatiotemporal abundance dynamics of adult delta smelt, we fit a suite of models to an extensive,

22    spatially resolved, catch survey time series from 13 annual cohorts. The best model included

23    cohort-specific abundance indicators and daily mortality rates, a regional spatial adjustment, and

24    haul-specific environmental conditions. The regional adjustment identified several density

25    hotspots that were persistent across cohorts. While this model did include local environmental

26    conditions, the gain in explained variation was relatively slight compared to that explained by the

27    regional adjustment. Total abundance estimates were derived by multiplying habitat volume by

28    catch density (design-based) and modeled density (model-based), with both showing severe

29    declines in the population over the time period studied. The design-based approaches had lower

30    uncertainty but potentially higher bias. We discuss the implications of our results for advancing

31    the science and improving management of delta smelt, and future data collection needs.

32

33    KEY WORDS: delta smelt; geostatistical models; population ecology; soap film smoothers; San

34    Francisco Estuary

## INTRODUCTION

Determining how and why an organism's population is distributed in space and time is a
fundamental organizing problem in population ecology (Krebs 1994). For small pelagic species
in tidal river estuaries, drawing inference about their distribution and abundance is especially
challenging because they cannot be tagged and tend to aggregate in schools that are small
relative to the size of their open-water habitats. Tidal river estuaries are ecotones characterized
by almost continuous multi-scale changes in environmental factors, from tidal to annual time
scales and with spatial scale changes ranging from the perceptual field of the organism up to the
entire span of the estuary (Odum 1988). In general, we can expect to need to apply models that
can disentangle the relative effects of processes acting at different spatiotemporal scales.

Multi-scale environmental variability can be especially important for small resident pelagic
species (Peebles et al. 2007; Reum et al. 2011). For example, tidal currents can influence vertical
and horizontal distributions so that organisms can maintain or change geographic position within
the estuary (Kimmerer and McKinnon 1987; Bennett and Burau 2015). Also, pelagic species will
move in response to temperature, turbidity, salinity and prey density gradients, all of which can
directly influence vital rates (Peebles et al. 2007; Reum et al. 2011; Rose et al. 2013) and shape
estuary wide and regional distributions. A practical consequence for model-based analysis of
distribution and abundance is that care must be taken to appropriately match the spatiotemporal
resolution of the data and the model. Models must include factors, and allow for predictions,
across multiple spatial and temporal scales simultaneously in order to provide useful insight into
spatiotemporal variability in abundance.

58    The San Francisco Estuary (SFE, Fig. 1) is a tidal river estuary ecotone with habitat composition

59    and structure that changes at multiple temporal and spatial scales (Cloern and Jassby 2012). One

60    of the largest tidal river estuaries on the west coast of the Americas, the SFE provides habitat for

61    delta smelt (*Hypomesus transpacificus*), an endemic annual pelagic fish species that inhabits the

62    low salinity and freshwater portions of the estuary upstream of San Pablo Bay (Fig. 1).

63    Substantial declines in the cohort abundance of delta smelt during the 1980s and early 1990s led

64    to protection under both the U.S. and California Endangered Species Acts in 1993, and new fish

65    monitoring programs, including one for the adult life stage. In addition to these spatially and

66    temporally extensive fish surveys, measurements of several salient environmental metrics have

67    also been collected in the SFE.

68

69    Delta smelt habitat preferences are relatively well understood (Moyle et al. 2016). The species

70    distribution is constrained down-estuary by salinity while up-estuary a variety of life stage

71    specific factors operate, including landward extent of tides, water clarity, salinity, temperature,

72    and risk of entrainment into water diversions (Sweetnam 1999; Bennett 2005; Kimmerer 2008;

73    Kimmerer 2008; Nobriga et al. 2008; Feyrer et al. 2011). Nevertheless, more precisely

74    understanding the spatiotemporal changes (or lack of changes) in abundance within the broader

75    range limits has been a focal point of conservation discussions (Brown et al. 2009; Feyrer et al.

76    2011), highlighting the need for statistical analysis at finer spatial and temporal scales than has

77    been typically carried out.

78

79    Our primary motivation was to gain insight into patterns of the distribution and abundance of

80    adult delta smelt. Specifically, we addressed three questions: Where do adult delta smelt

81    distribute themselves during their spawning season, and how variable is this distribution across

82    time (both within and between cohorts)? What factors operating at what scales most strongly

83    influence the spatial distributions? What are the year-over-year population growth rates?

84

85    To answer these questions, we constructed spatiotemporal models of catch density with three

86    different levels of spatiotemporal scale that we label global, regional, and local. Year of the

87    survey and cohort-specific mortality rates were global-level (i.e. population wide) components to

88    the model that described overall cohort specific trends in time. Regional (approximately 5 km

89    and larger) spatial variation is apparent from exploratory data analyses (Fig. 2) of catch per unit

90    volume (CPUV), the sum of all fish caught at given survey location divided by the sum of all

91    water sampled at that location, and this spatial variation was modeled using nonparametric

92    techniques. The importance of both within- and across-cohort changes in the regional spatial

93    distribution patterns were tested. At the local (individual sample) level, we estimated how much

94    of the variability in catch density was explained by three environmental covariates: water clarity,

95    salinity, and tide. Increased turbidity and decreased salinity are expected to have positive effects

96    on catch density based on *in situ* studies of earlier life stages (Nobriga et al. 2008; Feyrer et al.

97    2011). Flood relative to ebb tide was expected to increase catch densities as fish position

98    themselves within the water column and channel to either move upstream or to otherwise

99    maintain position (Feyrer et al. 2013; Bennett and Burau 2015). We also compared design-based

100   and model-based estimates of abundance for February of each year (definitions of design-based

101   and model-based inference are given in the supplementary material [SM] Section 4). Here the

102   aim was to quantify inter-annual changes and long-term trends, to assess how different

103 abundance estimates would be when standardizing effort for tide and to evaluate whether the two

104 approaches have any qualitative differences.

105

106 **METHODS**

107 *Data*

108 The California Department of Fish and Wildlife established the Spring Kodiak Trawl (SKT) in

109 2002 to collect data on the distribution and reproductive stage of spawning delta smelt. The SKT

110 survey usually visits 40 locations monthly from January through May (Fig. 1) over a several day

111 period. During each location visit a 10 minute surface trawl of the approximately top 2m of

112 water is taken. Three quarters of all trawls were made before noon. All delta smelt retained by

113 the gear are counted and measured, and the volume of water sampled (*Vol*, m$^3$) is estimated. We

114 used data from 2002-2014. Of the 2396 records used here, 1706 (71%) had zero catch. Of the

115 690 samples with positive catch, 227 recorded a single adult delta smelt, with a maximum catch

116 of 375.

117

118 The local tow-specific environmental covariate data included Secchi disk depth (*Sec*, cm), a

119 proxy for water clarity; specific conductance (*Cond*, microSiemens per centimeter, $\mu$S cm-1), a

120 proxy for salinity; and tide stage (*Tide*) which is categorically recorded as ebb (1500

121 observations), low slack (68 observations), high slack (97 observations), flood (731

122 observations). Although water temperature is also recorded, for this analysis we did not include it

123 in the models because the range of observed temperatures (min=6.6°C, max=23.6°C,

124 mean=12.9°C) were well within the tolerance of spawning and post-spawn adult delta smelt

125 (Swanson et al. 2000; Komoroske et al. 2014). Earlier versions of the model that did include

6

126  temperature never identified it as statistically significant. In contrast, measures of salinity up to

127  21ppt, high enough to constrain distribution and affect survival (Komoroske et al. 2014; Lisa M.

128  Komoroske et al. 2016), have been recorded in the SKT survey.

129

130  *Spatiotemporal catch density models*

131  The catch $y_{t,c,l}$ on Julian day $t$ of cohort $c$ at location $l$ was modeled using a negative binomial

132  distribution $y_{t,c,l} \sim \text{NegBin}(\mu_{t,c,l}, \theta)$ parameterized to have expected value $\mu$ and variance $\mu + \mu^2/\theta$

133  (Venables and Ripley 2002). The negative binomial was selected given evidence for

134  overdispersion relative to a Poisson distribution and from model residual diagnostics. The

135  different models for $\mu_{t,c,l}$ are described next and summarized in Table 1.

136

137  Most generally, the expected catches $\mu_{t,c,l}$ were modeled using a semi-parametric,

138  spatiotemporally explicit model within a generalized additive model (GAM) framework (Hastie

139  and Tibshirani 1986; Wood 2006; Augustin et al. 2013). The expected catch is the product of the

140  volume of water sampled, $Vol_{t,c,l}$, the true density $\delta_{t,c,l}$ in a spatially local region around $l$, and the

141  catchability $q_{t,c,l}$,

142  $$\mu_{t,c,l} = q_{t,c,l} \delta_{t,c,l} Vol_{t,c,l}. \quad (1)$$

143  Catchability $q_{t,c,l}$ has recently (Maunder et al. 2014) been conceptualized as a function of

144  availability (i.e. whether fish are in the tow path in the first place) and contact selectivity (the

145  probability that the net will catch and retain the fish given availability) (see Arreguín-Sánchez

146  1996 for other classic definitions). The catchability parameter $q_{t,c,l}$ is confounded with the overall

147  density, so it is assumed equal to1 for all the models. Further discussion of $q_{t,c,l}$ in the context of

148  adult delta smelt surveying is provided in the Discussion.

7

149

150    Modifications to Eqn. (1) were made to study different sources of variability in $\delta_{t,c,l}$. The first,

151    which is labeled global scale effects, and was included in all models and intended to capture

152    temporal trends in the overall density (total fish over total water volume), was to rewrite Eqn. (1)

153    as

$$\mu_{t,c,l} = \delta_{0,c}\exp(\beta_c t)Vol_{t,c,l} \qquad (2)$$

155    Eqn. (2) describes an exponential decline (assuming $\beta_c < 0$) in density from an overall initial

156    density $\delta_{0,c}$, and the expected catch is simply this density times the volume sampled on a given

157    tow.

158

159    An extension of the global density model of Eqn. (2) was to add a regional scale factor, namely a

160    dependency on space to the predictions,

$$\mu_{t,c,l} = \delta_{0,c}\exp(\beta_c t + s_{m,c})Vol_{t,c,l} \qquad (3)$$

162    where $s_{m,c} = s_{m,c}(UTMX_l, UTMY_l)$ is a nonparametric spatial smooth. A total of four different

163    hypotheses about how $s_{m,c}$ changed through time were considered: (1) it did not change in time,

164    $s_{m,c} = s$; (2) it depended only on the month of the survey, $s_{m,c} = s_m$; (3) it depended only on the

165    year of the survey, $s_{m,c} = s_c$; and (4) it depended on both the month and the year of the survey.

166    Because the spatial adjustments to the density vary at scales larger than the water surveyed in a

167    single trawl, these adjustments can be thought of as capturing spatially regional changes in

168    density.

169

170    The global and regional effects model given by Eqn. (3) was further extended to include local

171    scale effects. For each assumption about $s_{m,c}$, the effects of local environmental conditions on

172    $\delta_{t,c,l}$ were estimated with the model

173    $$\mu_{t,c,l} = \delta_{0,c}\exp(\beta_c t + s_{m,c} + \beta_{Sec}Sec_{t,c,l} + \beta_{Cond}Cond_{t,c,l} + \beta_{Tide(t,c),t,c,l}Tide_{t,c,l})Vol_{t,c,l}. \qquad (4)$$

174    The importance of Secchi and conductivity was also considered in the absence of a regional

175    spatial adjustment component, i.e. extending Eqn. (2) with these covariates.

176

177    In total fifteen different models were fit and evaluated (Table 1). Model fitting was done in the R

178    environment (R Core Team 2016) primarily using the glm.nb (Venables and Ripley 2002) and

179    gam (Wood 2004; Wood 2011) functions. Other functions and packages used are documented in

180    the model code provided in the SM. Soap film smoothers (Wood 2008) were used to make

181    spatial smooths $s_{m,c}$ follow large-scale habitat boundary features (SM Fig. S1). The boundaries

182    were set up in particular to avoid an influence of catch between Montezuma Slough and either

183    Cache Slough or Suisun Bay. Smoothing parameter estimation was done using maximum

184    likelihood (Wood 2011), but other criteria used for estimating the smooth parameter such as

185    generalized cross-validation did not qualitatively change the results. Secchi and conductivity

186    measurements were standardized to their z-scores prior to model fitting. A wide range of smooth

187    basis dimensions were considered to ensure results were not predicated on this choice, and

188    standard model residual diagnostics were investigated, including semivariograms (Clark 2007) of

189    residuals by month and year. Model comparison was done by assessing residual diagnostics,

190    Akaike's information criterion AIC (Burnham and Anderson 2002), fitted negative log-marginal-

191    likelihoods (NLML, see Eqn. 5 in Wood 2011).

192

193    Model evaluation of the effects of the locally measured covariates Secchi and conductivity was

194    partly complicated because of their global spatial structure. On average, more easterly (upstream)

195    regions of the delta smelt habitat are clearer and less saline (SM Fig. S2), leading to the

196    possibility that local environmental covariates will be confounded with the spatial terms in the

197    model. To approximate an upper bound on the most variability that local environmental

198    conditions might explain in the absence of spatial terms in the model, we computed the

199    proportion of null deviance explained by models of the form of Eqn. (2) but including each of

200    these covariates one at a time (Table 1 models 2-4). The proportion of the deviance explained by

201    each locally measured covariate when fitting the full model in Eqn. (4) (Table 1 models 13-15)

202    was also calculated by dropping each term individually and refitting the model while fixing the

203    smoothing parameters at the values estimated in model 9. This helped ensure that no changes in

204    the smoothing penalty upon refitting resulted in a "mopping up" of variation previously

205    accounted for by the removed covariate, thereby diminishing the estimated proportion of

206    deviance explained by the dropped covariate under consideration.

207

208    *Abundance estimation*

209    Total abundance estimates for the month of February for each year were made using both design-

210    based and model-based approaches (SM Sec. 4). Both approaches rely on volumetric expansions

211    of density estimates. The volumes were calculated by multiplying the area of water with at least

212    2 meters depth (provided by the United States Geological Survey) by 2 to compute the volume of

213    habitat $V_{tot}$ over which the density estimates might reasonably be extrapolated. This volume

214    excludes water deeper than 2 meters as well as shoal habitats. Thus our estimates are likely

215    underestimating the total population size depending on unknown densities in these unsampled

216    water volumes. However, this approach avoids extrapolating catch density information into

217    habitats that are not sampled by the SKT survey.

218

219    The design-based approach stratified the waterways most commonly occupied by delta smelt into

220    27 subregions (SM Fig. S3). The subregion, year and month specific catch densities were

221    expanded by subregion-specific water volumes and summed to obtain year and month-specific

222    abundance estimates. Assuming the abundance estimates were lognormally distributed, the 2.5

223    and 97.5 percentiles of this distribution were used to construct design-based prediction intervals.

224    Section 4.1 of the SM provides details on obtaining the parameters for these cohort specific

225    distributions.

226

227    In contrast to the design-based approach, the model-based approach does not require spatial

228    stratification of the habitat and allows predictions to be contingent on specific environmental

229    conditions thought to affect catchability. Based on model selection results, model 9 was used to

230    make model-based total abundance estimates as follows. We used 984 points distributed within

231    the spatial limits of the survey and the areas of water with at least 2 meters depth (SM. Fig S1) as

232    the spatial locations for predictions. At each one of these locations, the density per $10000m^3$ of

233    water was predicted on February $15^{th}$ (specifying a day is necessary for the Julian day effect) of

234    each year, the tide set equal to the flood factor level, and the Secchi and conductivity values

235    fixed at a month, year, and location specific value (described below). These densities were

236    averaged within each subregion, multiplied by the subregion water volume down to 2m, and

237    summed to produce overall abundance estimates (see SM Sec. 4.2 for details). Because direct

238    observations on Secchi depth and conductivity at the point locations used in making predictions

239    were not always available, spatially smoothed GAMs were used to predict both of these variables

240    during the February survey periods of each year. The GAMs were fit using the available survey

241    data on Secchi depth and conductivity and had the form $y_{t,c,l} \sim \text{Normal}(\mu = \beta_{m,c} + s_{m,c}, \sigma)$, where

242    $y_{t,c,l}$ was either the z-transformed Secchi depth or conductivity measurements from the SKT

243    survey. The fits were generally quite good: the models of Secchi depth and conductivity

244    described at least 88% and 94% of the null deviance for 80% of the months, respectively.

245    Abundance prediction intervals were estimated using a parametric bootstrapping approach that

246    included uncertainty in model parameters, covariate predictions, and observations (see SM Sec.

247    4.2).

248

249    **RESULTS**

250    Table 1 shows model summary statistics. There was clear support for including both a regional

251    spatial adjustment and local environmental conditions in the expected catch models. The best

252    model identified by AIC included a separate spatial distribution for each month (model 10),

253    while the negative log-marginal-likelihood identified a model with a constant spatial smooth

254    over time as the best (model 9). Residual diagnostics for models without a regional smooth

255    adjustment term were poor as measured by distributional checks of residuals. In contrast, models

256    including a regional spatial term had residual qq-plots and semivariograms that suggested no

257    systematic bias in predictions due to the spatial variability in the distribution. Simpler models

258    had higher dispersion parameters, reflecting larger prediction error when the mean structure was

259    less flexible.

260

261 Models including a smooth term to capture regional variation in catch identified several density

262 hotspots (Figs. 3 and S4; see also Fig. 2 for empirical densities): the waterways surrounding

263 Grizzly Island, channels at the confluence of the Sacramento and San Joaquin rivers, the Cache

264 Slough complex, and the Sacramento deep water shipping channel. These density hotspots were

265 fairly consistent between cohorts, with the Cache Slough complex and Sacramento deep water

266 shipping channel the most persistently high. We focused on model 9 for making predictions

267 because the differences in month-specific predictions in model 10 are dominated by

268 disappearance of density hotspots in April and May (likely reflecting post-spawning mortality)

269 rather than a spatial shift in the locations of hotspots (SM Fig. S4).

270

271 The local environmental covariates tide and conductivity explained very little (<2%) of the null

272 deviance beyond that of model 1, but Secchi depth explained an additional 21.3% of the null

273 deviance when no regional spatial adjustment was made (Table 1, models 2-4 and 13-15). The

274 effect size on the linear scale of Secchi was approximately double that of conductivity, but both

275 local covariates could translate into substantially larger expected changes in density predictions

276 over the range of observed turbidity and salinity indices (Table 2). Catch density was higher on

277 flood and low slack tide levels in comparison with ebb tide (the increase on low slack tide was

278 the highest, but surveys during this tide stage account for <3% of samples), and not significantly

279 different for high slack tide conditions (Table 2).

280

281 Figure 4 shows the total abundance estimates and prediction uncertainty for February 15[th] of

282 each year (see SM Table S1 for values) for the design- and model-based estimates. The

283 geometric mean annual growth rate over the 13 years was 0.88 and 0.87 for the design- and

13

284    model-based approaches, respectively, and the percentage decline from 2002 to 2014 was 82%

285    and 79% for the design- and model-based approaches, respectively. Note that the results about

286    declines do not depend on the tide factor level choice used in making total abundance estimates.

287    Despite the general agreement between design- and model-based estimates of trend, the two

288    approaches showed the same annual growth rate in only 6 of the 13 years, and differed in

289    magnitude especially in 2003 and 2012 (Fig. 4 and SM Table 1). The differences in abundance

290    magnitude did depend on the model chosen, with the most complicated model showing

291    predictions very similar to the design-based approach (SM Fig. S5).

292

293    **DISCUSSION**

294    For small, elusive, and rare pelagic fish species such as delta smelt, often the only source of

295    information from the wild is catch density from trawls or other types of nets (e.g., beach seining),

296    along with additional measurements of local environmental conditions. Given such data, at a

297    minimum we would like to quantify the variability in distribution and abundance. Ideally, we

298    could go further to identify causal factors that explain the variability at different scales, or rule

299    out those that do not, and to assess the extent to which findings from theoretical and laboratory

300    work are identifiable in the wild.

301

302    The spatial distributions quantified here are similar to the descriptive reports by Merz et al.

303    (2011) and Murphy and Hamilton (2013) in their general depiction. By constructing statistical

304    models, we were able to test hypotheses about the variability of this spatial structure. At a

305    regional scale, our models indicated that the distribution of adult delta smelt was fairly consistent

306    across months and years, with the dominant within-year change being disappearance of hotspots

14

307    likely due to post-spawn mortality as the spawning season progresses for this annual species.

308    This suggests that the majority of regional movement from juvenile and sub-adult rearing

309    locations to spawning areas has already happened by the time the SKT survey is conducted, that

310    spawning habitat locations are relatively constant within and between years, and that no

311    substantial further restructuring of the population at regional scales occurs afterwards.

312

313    What leads to the emergence of density hotspots remains to be determined. A recent pairing of

314    the sub-adult delta smelt catch data used by Feyrer et al. (2011) with a three-dimensional

315    hydrodynamic model suggests that density hotspots may reflect the interplay of local water

316    quality conditions with tidal velocity differences that exist between shoals and deeper shipping

317    channels (Bever et al. 2016). Other possible explanations for adult and spawning delta smelt

318    spatial variation include distributions of prey or spawning habitats, or areas more suited for

319    survival during spawning. Why no density hotspots emerge and persist upstream of the Jersey

320    Point (located near the arrow tip showing the San Joaquin River in Fig. 1) area remains to be

321    determined, but likely factors include inhospitable habitat and advection of fish into water export

322    facilities (Kimmerer 2008; Kimmerer 2011).

323

324    At local spatial scales there continues to be high variability in the spatial distribution (which

325    necessitated the use of a negative binomial catch distribution model), some of which is likely

326    related to spawning-related aggregations of delta smelt and some of which is related to changes

327    in local salinity (movement away from) and turbidity (movement towards) conditions. Our view

328    is that the best interpretation of the categorical covariate tide is that it affects changes in fish

329    availability to the gear, a component of catchability $q_{t,c,l}$, with the direction of the effects found

330    here being consistent with Feyrer et al. (2013). In general it appears that, due to its relatively

331    coarse spatial and temporal resolution, the SKT survey cannot distinguish between very local,

332    site level movement, up to movement between adjacent locations, and changes in catchability

333    related to local environmental conditions. The infrequent yet extremely large catches point to

334    highly localized and ephemeral aggregations of fish but, similar to questions about the existence

335    of regional density hotspots, the relative contributions of social cues vs. habitat cues vs.

336    hydrodynamics leading to the formation of these aggregations remains to be determined.

337

338    Previous analyses of the sub-adult life-stage have found local environmental covariates to be

339    statistically significant predictors of delta smelt distribution, with Feyrer et al. (2011) remarking

340    that "specific conductance and Secchi depth accounted for a meaningful reduction of null

341    deviance." In contrast, we found that these covariates explained very little of the variation in

342    adult catch when a regional spatial adjustment to density was included. The comparatively large

343    amount of deviance explained by Secchi depth when no spatial smooths were included in the

344    model (model 2) suggests that water clarity has some influence on both local and regional

345    distributions, although from a statistical perspective any models not containing a spatial

346    adjustment beyond what is made by the local environmental covariates were very inferior. While

347    suitable local environmental conditions are necessary to explain the distribution and abundance

348    of delta smelt, they are far from sufficient. We suggest that to better understand both the regional

349    and local changes in densities, an understanding of the characteristics leading to ideal spawning

350    habitat features is needed, along with assessments of the variability of these characteristics in

351    space and time.

352

353     At the decadal time scale delta smelt are currently in a severe state of population decline, with

354     suspected causes including removal of water from the system and alien species (Moyle et al.

355     2016). Here we used the best available survey data to quantify this decline more precisely.

356     Design- and model-based approaches closely agreed in the rate and amount of overall decline

357     from 2003 to 2014.

358

359     Despite the general agreement in long-term trends between the two approaches for abundance

360     estimation, there were also differences. In 2003 the design-based estimates showed a decline in

361     abundance compared to 2002, while model-based estimates showed an increase. During this year

362     the frequency of sampling on the flood tide was only 8%, and this may have led to the qualitative

363     mismatch in year-over-year abundance change between the design- and model-based methods. It

364     seems likely that the design-based approach is negatively biased when compared with the model-

365     based approaches due to the failure to account for the effect of tide cycle on catchability $q_{t,c,l}$.

366     Another difference was in prediction intervals, with model-based ones being notably wider likely

367     related to the more complete inclusion of the different sources of uncertainty in the model-based

368     approach which is accounting for spatial and tow-specific sources of uncertainty. Finally, the

369     magnitude of the estimates also differed, with model-based estimates generally being

370     substantially higher, although models with more complicated smooths had estimates that

371     increasingly approached the design-based ones (SM Fig. S5). This closer agreement of the

372     models with the most complicated smooths and the design-based approach is likely due in part to

373     overfitting, whereby the expected model predictions are able to more closely track zero catch

374     data. Other surveys making multiple tows per site visit have found that although the frequency of

375     zero catch was similarly high on any given tow, nonzero catch usually occurred at least once

17

376   (Polansky et al. 2014). Thus, we suspect that the models with simpler spatial smooth terms are

377   more reflective of actual distributions because they are drawing on information across time, and

378   hence less informed by zero catch when in fact fish may be locally in the area. Whether using

379   design- or model-based approaches to construct abundance estimates, information about false

380   zero catches as well as abundances in shoal habitats as well as the vertical density gradients in

381   channel and open-water habitats are needed to reduce abundance estimate bias and uncertainty.

382

383   Pinpointing the relative contributions of anthropogenic vs. natural sources to the population

384   decline will continue to be challenging, and will likely best be done in a complete life-cycle

385   analysis framework that integrates survey data from all life stages. Absolute abundance estimates

386   will first be needed from each source in order to integrate information from different life stages,

387   and catch level models such as applied here can help achieve this. The importance of tide, found

388   here and elsewhere (Bennett et al. 2002; Feyrer et al. 2013; Bennett and Burau 2015),

389   emphasizes a need to consider accounting for this covariate analyses where organism detection

390   might be driven by tidal conditions (see also Arreguín-Sánchez 1996) to control for its effect on

391   catch density. None of the previous population dynamics models using annual abundance indices

392   (Mac Nally et al. 2010; Thomson et al. 2010; Maunder and Deriso 2011) attempted to

393   standardize catch data when making these indices, which could mean that abundance and

394   covariate relationships have not been described accurately.

395

396   **Acknowledgments**

18

399 the authors and do not necessarily reflect the opinions of the U.S. Department of the Interior or

400 the U.S. Fish and Wildlife Service.

401

402 **Supplementary Material**

403 Supplementary material with additional details, figures and code is provided. Data and code are

404 available from the U.S. Fish and Wildlife Service.

405

406 **REFERENCES**

407

408 Arreguín-Sánchez, Francisco. 1996. Catchability: a key parameter for fish stock assessment.

409     *Reviews in Fish Biology and Fisheries* 6: 221–242.

410 Augustin, Nicole H., Verena M. Trenkel, Simon N. Wood, and Pascal Lorance. 2013. Space-time

411     modelling of blue ling for fisheries stock management. *Environmetrics* 24: 109–119.

412 Bennett, William A., and Jon R. Burau. 2015. Riders on the storm: selective tidal movements

413     facilitate the spawning migration of threatened delta smelt in the San Francisco Estuary.

414     *Estuaries and Coasts* 38: 826–835.

415 Bennett, William A. 2005. Critical assessment of the delta smelt population in the San Francisco

416     Estuary, California. *San Francisco Estuary and Watershed Science* 3: 1–71.

417 Bennett, William A., Wim J. Kimmerer, and Jon R. Burau. 2002. Plasticity in vertical migration

418     by native and exotic estuarine fishes in a dynamic low-salinity zone. *Limnology and*

419     *Oceanography* 47: 1496–1507.

420 Bever, Aaron J., Michael L. MacWilliams, Bruce Herbold, Larry R. Brown, and Frederick V.

421     Feyrer. 2016. Linking hydrodynamic complexity to delta smelt (*Hypomesus*

422      *transpacificus*) distribution in the San Francisco Estuary, USA. *San Francisco Estuary*

423      *and Watershed Science* 14: 1–25.

424      Brown, Larry R., Wim J. Kimmerer, and Randall Brown. 2009. Managing water to protect fish: a

425      review of California's environmental water account, 2001-2005. *Environmental*

426      *Management* 43: 357–368.

427      Burnham, Kenneth P., and David Anderson. 2002. *Model selection and multimodel inference: a*

428      *practical information-theoretic approach*. New York, NY: Springer.

429      Clark, James S. 2007. Models for Ecological Data: An Introduction. Princeton, NJ: Princeton

430      University Press.

431      Cloern, James E., and Alan D. Jassby. 2012. Drivers of change in estuarine-coastal ecosystems:

432      Discoveries from four decades of study in San Francisco Bay. *Reviews of Geophysics* 50.

433      Feyrer, Frederick, Ken B. Newman, Matthew Nobriga, and Ted Sommer. 2011. Modeling the

434      effects of future outflow on the abiotic habitat of an imperiled estuarine fish. *Estuaries*

435      *and Coasts* 34: 120–128.

436      Feyrer, Frederick, Donald Portz, Darren Odum, Ken B. Newman, Ted Sommer, Dave Contreras,

437      Randall Baxter, Steven B. Slater, Deanna Sereno, and Erwin Van Nieuwenhuyse. 2013.

438      SmeltCam: underwater video codend for trawled nets with an application to the

439      distribution of the imperiled delta smelt. *PLoS ONE* 8: e67829.

440      Hastie, Trevor, and Robert Tibshirani. 1986. Generalized additive models (with discussion).

441      *Statistical Science* 1: 297–318.

442      Kimmerer, Wim J. 2008. Losses of Sacramento River Chinook Salmon and Delta Smelt to

443      Entrainment in Water Diversions in the Sacramento–San Joaquin Delta. *San Francisco*

444      *Estuary and Watershed Science* 6.

445      Kimmerer, Wim J. 2011. Modeling Delta Smelt Losses at the South Delta Export Facilities. *San*

446          *Francisco Estuary and Watershed Science* 9.

447      Kimmerer, Wim J., and A. D. McKinnon. 1987. Zooplankton in a marine bay. II. Vertical

448          migration to maintain horizontal distributions. *Marine Ecology Progress Series* 41: 53–

449          60.

450      Komoroske, L. M., R. E. Connon, J. Lindberg, B. S. Cheng, G. Castillo, M. Hasenbein, and N.

451          A. Fangue. 2014. Ontogeny influences sensitivity to climate change stressors in an

452          endangered fish. *Conservation Physiology* 2:1-13.

453      Komoroske, Lisa M., Ken M. Jeffries, Richard E. Connon, Jason Dexter, Matthias Hasenbein,

454          Christine Verhille, and Nann A. Fangue. 2016. Sublethal salinity stress contributes to

455          habitat limitation in an endangered estuarine fish. *Evolutionary Applications* 9: 963–981.

456      Krebs, Charles J. 1994. *Ecology: the experimental analysis of distribution and abundance*. 4th

457          ed. New York, NY: HarperCollins College Publishers.

458      Mac Nally, Ralph, James R. Thomson, Wim J. Kimmerer, Frederick Feyrer, Ken B. Newman,

459          Andy Sih, William A. Bennett, et al. 2010. Analysis of pelagic species decline in the

460          upper San Francisco Estuary using multivariate autoregressive modeling (MAR).

461          *Ecological Applications* 20: 1417–1430.

462      Maunder, Mark N., and Richard B. Deriso. 2011. A state–space multistage life cycle model to

463          evaluate population impacts in the presence of density dependence: illustrated with

464          application to delta smelt (*Hyposmesus transpacificus*). *Canadian Journal of Fisheries*

465          *and Aquatic Sciences* 68: 1285–1306.

466    Maunder, M.N., P.R. Crone, J.L. Valero, and B.X. Semmens. 2014. Selectivity: Theory,

467        estimation, and application in fishery stock assessment models. *Fisheries Research* 158:

468        1–4.

469    Merz, Joseph E., Scott Hamilton, Paul S. Bergman, and Bradley Cavallo. 2011. Spatial

470        perspective for delta smelt: a summary of contemporary survey data. *California Fish and*

471        *Game* 97: 164–189.

472    Moyle, Peter B., Larry R. Brown, John R. Durand, and James A. Hobbs. 2016. Delta smelt: life

473        history and decline of a once-abundant species in the San Francisco Estuary. *San*

474        *Francisco Estuary and Watershed Science* 14: 1–30.

475    Murphy, Dennis D., and Scott A. Hamilton. 2013. Eastward migration or marshward dispersal:

476        exercising survey data to elicit an understanding of seasonal movement of delta smelt.

477        *San Francisco Estuary and Watershed Science* 11: 1–21.

478    Nobriga, Matthew L., Ted R. Sommer, Frederick Feyrer, and Kevin Fleming. 2008. Long-term

479        trends in summertime habitat suitability for delta smelt (*Hypomesus transpacificus*). *San*

480        *Francisco Estuary and Watershed Science* 6: 1–13.

481    Odum, William E. 1988. Comparative Ecology of Tidal Freshwater and Salt Marshes. *Annual*

482        *Review of Ecology and Systematics* 19: 147–176.

483    Peebles, Ernst B., Scott E. Burghart, and David J. Hollander. 2007. Causes of Interestuarine

484        Variability in Bay Anchovy (*Anchoa mitchilli*) Salinity at Capture. *Estuaries and Coasts*

485        30: 1060–1074.

486    Polansky, L., Matt Nobriga, Ken Newman, Matt Dekar, Kim Webb, and Mike Chotkowski.

487        2014. Delta smelt movement during and extreme drought: Intensive Kodiak trawling

488        at Jersey Point. *Interagency Ecological Newsletter* 4:5-13.

489     R Core Team. 2016. *R: A Language and Environment for Statistical Computing* (version 3.3.0).

490             Vienna, Austria: R Foundation for Statistical Computing.

491     Reum, Jonathan C. P., Timothy E. Essington, Correigh M. Greene, Casimir A. Rice, and Kurt L.

492             Fresh. 2011. Multiscale influence of climate on estuarine populations of forage fish: the

493             role of coastal upwelling, freshwater flow and temperature. *Marine Ecology Progress*

494             *Series* 425: 203–215.

495     Rose, Kenneth A., Wim J. Kimmerer, Karen P. Edwards, and William A. Bennett. 2013.

496             Individual-based modeling of delta smelt population dynamics in the upper San Francisco

497             Estuary: I. model description and baseline results. *Transactions of the American*

498             *Fisheries Society* 142: 1238–1259.

499     Swanson, Christina, Turid Reid, Patricia S. Young, and Joseph J. Cech Jr. 2000. Comparative

500             environmental tolerances of threatened delta smelt (*Hypomesus transpacificus*) and

501             introduced wakasagi (*H. nipponensis*) in an altered California estuary. *Oecologia* 123:

502             384–390.

503     Sweetnam, Dale A. 1999. Status of delta smelt in the Sacramento-San Joaquin Estuary 85: 22–

504             27.

505     Thomson, James R., Wim J. Kimmerer, Larry R. Brown, Ken B. Newman, Ralph Mac Nally,

506             William A. Bennett, Frederick Feyrer, and Erica Fleishman. 2010. Bayesian change point

507             analysis of abundance trends for pelagic fishes in the upper San Francisco Estuary.

508             *Ecological Applications* 20: 1431–1448.

509     Venables, W.N., and B.D. Ripley. 2002. Modern applied statistics with S. Fourth Edition. New

510             York, NY: Springer.

511  Wood, Simon N. 2004. Stable and efficient multiple smoothing parameter estimation for

512      generalized additive models. *Journal of the American Statistical Association* 99: 673–

513      686.

514  Wood, Simon N. 2006. Generalized additive models: an introduction with R. Boca Raton, FL:

515      Chapman & Hall.

516  Wood, Simon N. 2008. Soap film smoothing. *Journal of the Royal Statistical Society B* 70: 931–

517      955.

518  Wood, Simon N. 2011. Fast stable restricted maximum likelihood and marginal likelihood

519      estimation of semiparametric generalized linear models. *Journal of the Royal Society:*

520      *Series B* 73: 3–36.

**Tables**

522  Table 1- Overview of catch models showing the (effective) degrees of freedom (df), information theoretic measures of model

523  goodness of fit (AIC and ΔAIC), the negative log-marginal-likelihood (NLML- smaller values are better), and percent of the null

524  deviance explained (% dev. expl.). Global is defined by Eqn. (2) and Global + regional is defined by Eqn. (3), while local terms are

525  Secchi (Sec), conductivity (Cond), and Tide- see Eqn. (4). Regional spatial smooth terms were either constant across months and years

526  (single), different by month but not year (monthly), different by year but not month (yearly), or different for each month and year.

527  Models 13-15 had fixed smooth term parameters using the estimates from model 9 and were used for estimating the % dev. expl. by

528  each of the three individual local terms in model 9.

| Model | Density model | df | AIC | ΔAIC | NLML | % dev. expl. | $\theta$ |
|---|---|---|---|---|---|---|---|
| 1 | Global | 27 | 6717.2 | 1178.2 | - | 12.9 | 0.1 |
| 2 | Global + Sec | 28 | 6337.8 | 798.9 | - | 34.2 | 0.2 |
| 3 | Global + Cond | 28 | 6692.8 | 1153.8 | - | 14.5 | 0.1 |
| 4 | Global + Tide | 30 | 6701.4 | 1162.5 | - | 14.2 | 0.1 |
| 5 | Global + regional (single) | 49.9 | 5643.0 | 104.0 | 2821.8 | 63.6 | 0.4 |
| 6 | Global + regional (monthly) | 118.1 | 5638.3 | 99.3 | 2853.3 | 67.4 | 0.4 |
| 7 | Global + regional (yearly) | 199.0 | 5603.3 | 64.4 | 2831.6 | 72.2 | 0.5 |
| 8 | Global + regional (month and year) | 632.5 | 5888.7 | 349.7 | 2933.1 | 83.8 | 0.8 |
| 9 | Global + regional (single) + Sec + Cond + Tide | 54.7 | 5548.2 | 9.3 | 2769.2 | 66.9 | 0.4 |
| 10 | Global + regional (monthly) + Sec + Cond + Tide | 128.1 | 5538.9 | 0.0 | 2789.3 | 70.6 | 0.5 |
| 11 | Global + regional (yearly) + Sec + Cond + Tide | 198.6 | 5572.3 | 33.4 | 2798.6 | 72.9 | 0.5 |
| 12 | Global + regional (month and year) + Sec + Cond + Tide | 506.3 | 5726.6 | 187.7 | 2819.6 | 82.4 | 0.7 |
| 13 | Global + regional (single, fixed) + Cond + Tide | 52.7 | 5606.3 | 67.4 | 2801.1 | 65.0 | 0.4 |
| 14 | Global + regional (single, fixed) + Sec + Tide | 52.7 | 5566.9 | 28.0 | 2780.2 | 66.1 | 0.4 |
| 15 | Global + regional (single, fixed) + Sec + Cond | 50.7 | 5557.9 | 18.9 | 2778.0 | 66.4 | 0.4 |

529    Table 2- Parameter estimates and bootstrapped estimates of uncertainty for the parameters

530    associated with the local environmental covariates for model 9 (see Table 1) on the $\log_e$ scale.

531    Lower and upper columns show the 2.5 and 97.5 percentiles from 1000 samples from a

532    multivariate normal distribution parameterized by the mean and covariance matrix from the fitted

533    model 9. The final columns show density prediction differences on the response scale given the

534    described local environmental change, where the changes are based on changing from the 2.5 to

535    the 97.5 percentile for the continuous covariate observations, and in comparison with an ebb tide.

| Covariate | Estimate | Lower | Upper | Density factor change on response scale | |
|---|---|---|---|---|---|
| Secchi depth | -0.880 | -1.112 | -0.670 | Decrease in turbidity | 0.415 |
| Conductivity | -0.403 | -0.583 | -0.232 | Increase in salinity | 0.669 |
| Flood | 0.338 | 0.113 | 0.552 | From ebb to flood | 1.398 |
| High slack | -0.093 | -0.658 | 0.476 | From ebb to high slack | 0.910 |
| Low slack | 0.962 | 0.389 | 1.571 | From ebb to low slack | 2.622 |

536

**Figures**

Figure 1- Overview of the inland portion of the San Francisco Estuary where adult delta smelt are most commonly found. Black x's denote the regular monthly Spring Kodiak Trawl survey locations.

Figure 2- Mean catch per unit volume at each sampling location for each month (averaged over 2002-2014). Units are per $10000m^3$ of water.

Figure 3- Density predictions at a flood tide per $10000m^3$ of water based on model 9 on February 15th 2004 using the mean Secchi and conductivity values. By fixing the local covariates the figure emphasizes density variation due to intrinsic variability. For clarity catch densities above 10 fish/$10000m^3$ of water are colored the same. See SM Fig. S4 for month specific predictions using model 10.

Figure 4- Abundance estimates on February 15th of each year. Design-based abundance estimates are shown by the line with filled circles with vertical lines extending to the 2.5 and 97.5 percentiles of the lognormal distributions. Model-based predictions from model 9 are shown as a solid line with dashed lines drawn at the 2.5 and 97.5 prediction percentiles based on 1000 bootstrapped samples. Inset numbers show the percentage of samples in each February that were done on a flood tide to illustrate the variability in sample conditions, which the model-based estimates account for. See SM Fig. S5 predictions using models 8 and 12.
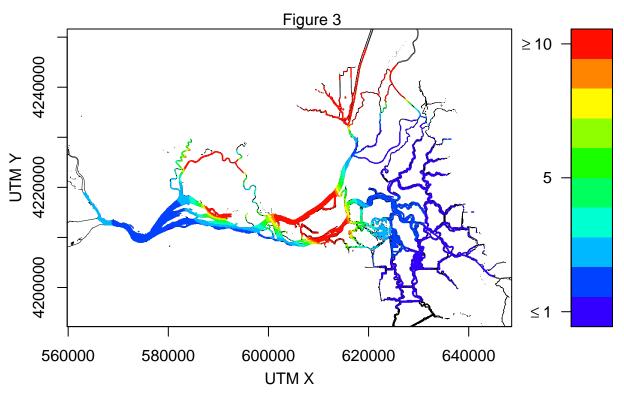
Figure 1



Figure 1

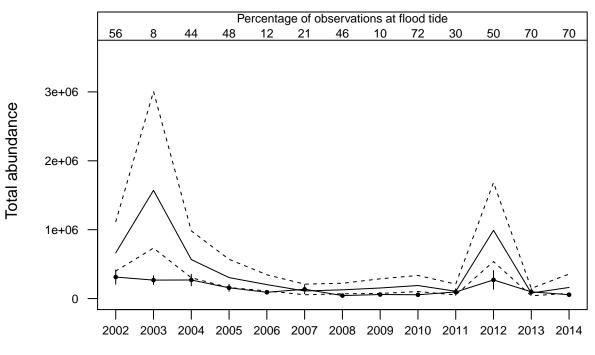Figure 2

# Figure 2
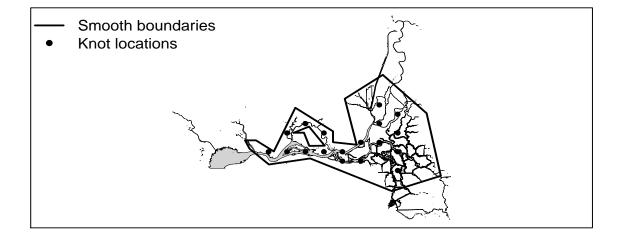
Figure 3



Figure 3

Figure 4

Figure 4

1   **Spatiotemporal models of an estuarine fish species to identify patterns and**
2   **factors impacting their distribution and abundance**

3   Leo Polansky, Ken B. Newman, Matthew L. Nobriga, Lara Mitchell

4   Supplementary Material

# 5 1  Survey locations, knots for smoothing basis, boundaries, and density pre-
6   diction locations

Figure S1: Model smooth boundaries, knot locations, and prediction locations used in the analysis.

Figure S2: Tow specific values of Secchi and electrical conductance vs. UTMX. Points are colored by the month during which they were recorded: purple-January; blue-February; cyan-March; April-green; May-orange.



## 2   Locally measured covariates

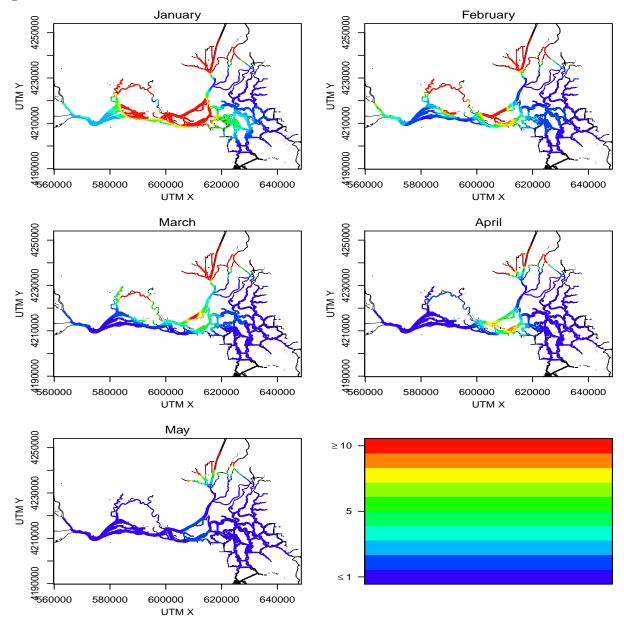A visual display of how turbidity and salinity vary in the UTMX direction, which corresponds approximately to an up and down estuary change, is shown in Figure S2.

## 3   Intra-cohort distribution changes

Figure S3 shows density predictions using model 10 of the main text, which has a different spatial smooth for each month.

## 4   Total abundance estimates

Two distinctly different perspectives in sampling theory on making inferences from samples to populations are design-based inference and model-based inference (Thompson 2002). Design-based inference views the values on sampling units as fixed, non-random quantities, and the only randomness present is that induced by the sample selection pro-

2

Figure S3: Month specific density predictions based on model 10 (Table 1 of the main text) at a flood tide per 10000m$^3$ using the mean Secchi and conductivity values; compare with figure 3 of the main text. Turbidity and salinity at each prediction point are set at their mean value, the Julian day is 15, 45, 74, 105, and 135 for the months of January, February, March, April, and May, respectively. As such these density maps emphasize the changes in density due to spatial and temporal changes. Color value map is shown in figure 3 of the main text.

18 cess. For example, assume in a body of water there are N cubic meter water volumes
19 "plots", with plot $i$ having $y_i$ individual fish, $i = 1, ..., N$, and the inference objective is
20 to estimate total population, $\tau_y = \sum_{i=1}^{N} y_i$. A simple random sample of size $n < N$ is
21 drawn without replacement and $\tau_y$ is estimated by multiplying the sample average of $y$ by
22 $N$. The total population estimate is thus a random variable where the randomness arises
23 solely from the random selection process.

24 In contrast, model-based inference views the values on sample units as realizations from
25 some underlying random natural process. When the sample units are partitions of a spatial
26 domain the random process often induces spatial correlation in the attributes defined on
27 the units, e.g., adjacent plots are more likely to have similar values than more spatially
28 separated plots. Inference is directed at estimating parameters that characterize the
29 underlying random natural process, e.g., a mean value ($\mu$), variance ($\sigma^2$), and covariance
30 between plots $i$ and $j$ ($\sigma_{i,j}$). Realized population characteristics, e.g., $\tau_y$, can still be
31 estimated using estimates of the parameters of the random process, e.g. $\hat{\tau} = N\hat{\mu}$.

32 We note that strictly speaking, from a model-based inference perspective, the sample
33 units do not need to be randomly selected for inference. However, it is our view that such
34 additional human-induced randomization is advisable as it allows for comparison between
35 model-based and design-based inference, and assessment of the sensitivity of assumptions
36 made about the random process.

## 37  4.1   Design-based total abundance estimates

Design-based estimates of total monthly abundance $N_{tot}$ (indices for month and year are
suppressed for clarity) were calculated with historical SKT data by dividing the delta
into 27 subregions (see Fig. S5) and carrying out volume expansions of average delta
smelt catch densities at the subregion level. The average density in each subregion was
calculated as the total catch divided by the total water volume sampled

$$\hat{\delta}_h = \frac{\sum_{j=1}^{m_h} Catch_{h,j}}{\sum_{j=1}^{m_h} Vol_{h,j}}$$

is the average density calculated over the $m_h$ sampling locations in the subregion, $Catch_{h,j}$
is the catch in a single tow $j$ in subregion $h$, and $Vol_{h,j}$ is the associated tow volume.
The total abundance was calculated by expanding the subregion specific catch densities
by the water volume in areas at least 2 meters deep down to 2 meters depth in subregion
$h$, $Vol_h$, and then summing across all subregions the subregion specific totals $\hat{N}_h$,

$$\hat{N}_{tot} = \sum_{h=1}^{27} \hat{N}_h = \sum_{h=1}^{27} \hat{\delta}_h Vol_h$$

38 In some months not all 27 subregions were sampled by the SKT. In cases where subregion
39 density estimates were missing due to lack of sampling, an estimate from a neighboring
40 subregion was used for imputation.

4

For a given year, variance estimates for the total abundances are given by

$$Var\left(\hat{N}_{tot}\right) = \sum_{h=1}^{h=29} \left( \frac{Vol_h^2 \, s_h^2}{\left(\frac{1}{m_h} \sum_{j=1}^{n_h} Vol_{h,j}\right)^2 m_h} \right)$$

where

$$s_h^2 = \frac{\sum_{j=1}^{m_h} \left(Catch_{h,j} - \hat{\delta}_h Vol_{h,j}\right)^2}{m_h - 1}$$

41 is the variance contribution from each subregion. Some values of $s_h^2$ were missing because
42 no sampling was done in a subregion or because a single site was sampled (in which case
43 $m_h = 1$). In these cases, the median value of $s_h^2$, calculated over all available values, was
44 used in place of missing values.

Suppressing time specific indices, confidence intervals were calculated for these abun-
dance estimates by assuming the abundances $\hat{N}_{tot}$ were log-normally distributed. For a
sample point estimate $\hat{N}_{tot}$ and variance of $Var\left(\hat{N}_{tot}\right)$, an estimate of the coefficient of
variation is

$$CV = \frac{\sqrt{Var(\hat{N}_{tot})}}{\hat{N}_{tot}}$$

the location and scale parameters are

$$\mu = \log_e \left(\frac{\hat{N}_{tot}}{\sqrt{1 + CV^2}}\right)$$

and

$$\sigma = \sqrt{\log_e \left(1 + CV^2\right)}$$

45 respectively. The natural log transformed abundance $z = \log_e(\hat{N}_{tot})$ is normally dis-
46 tributed with mean $\mu - \sigma^2/2$ and variance $\sigma^2$, where the mean is bias corrected so that
47 the expected value of $\exp(z)$ is $\hat{N}_{tot}$. The 2.5 and 97.5 percentiles of z were exponentiated
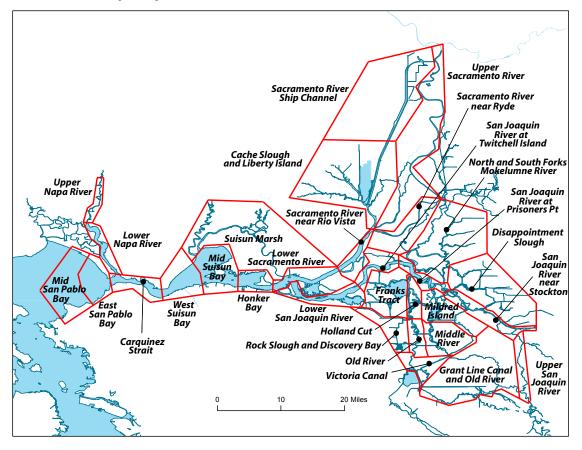48 to estimate a 95% confidence interval for $\hat{N}_{tot}$.

49 ## 4.2    Model-based total abundance estimates

Similar to the design-based approach using subregion specific average catch densities as
the starting point for constructing a total abundance estimate, the model-based approach
uses a model averaged density estimate per subregion to expand by subregion specific
water volumes en route to obtaining a total abundance estimate. The month and year
indices are suppressed for clarity. Denote the parameter vector of coefficients (including
the coefficients for the smooth terms) from the fitted model by

$$\hat{\boldsymbol{\beta}} = \left[\hat{\beta}_0, \ldots, \hat{\beta}_n\right]^\top$$

5

Figure S4: Spatial stratification of the Delta used for subregion based density expansions in the design-based estimates of abundance. The Mid and East San Pablo Bay subregions, along with Franks Tract, were excluded in total abundance calculations because these areas are not surveyed by the SKT.



and denote by $\mathbf{X}_i$ the design vector of values at location $i$ from one of the lcoations shown in Fig S1. The model estimated abundance on the log scale at location $i$ is

$$\hat{Y}_i = \mathbf{X}_i\hat{\boldsymbol{\beta}} + \log_e\left(Vol_p\right)$$

where $Vol_p$ is the prediction volume. The estimated number of fish $\hat{y}_i$ (per $Vol_p\mathrm{m}^3$ of water) at location $i$ on the response scale is

$$\hat{y}_i = \exp\left(\hat{Y}_i\right)$$

The mean density in subregion $h$ is

$$\hat{\bar{\delta}}_h = \frac{\sum_{k=1}^{K_h} \hat{y}_i}{K_h Vol_p}$$

where the sum is over the $K_h$ locations in subregion $h$. The estimate of the total abundance in subregion $h$ is

$$\hat{N}_h = \hat{\bar{\delta}}_h Vol_h$$

6

where $Vol_h$ is the total water volume in areas at least 2 meters deep down to 2 meters depth in subregion $h$. Again the total abundance estiamte is simply the sum of these subregion level estimates over all subregions

$$\hat{N}_{tot} = \sum_{h=1}^{27} \hat{N}_h$$

Prediction intervals for total abundance were obtained by parametric bootstrap and posterior simulation of GAM model coefficients, conditional on the smoothing parameter (Sections 4.8 and 5.4.2 in Wood (2006) and the scale parameters (the $\hat{\theta}$ from the catch model GAM and the $\hat{\sigma}$'s from the Secchi and conductivity GAMs). (Obtaining posterior distributions unconditional on the smoothing parameter, as outlined in Section 4.9.3 in Wood (2006) which involves wrapping the entire steps described next into a simulation-refitting process was not possible due to computational time.) Because estimates of total abundance are based on predicted values of Secchi depth and conductivity at each point location $i$, the first step is to simulate predictions of these covariate values at each location, and then, given these values, simulate location specific observations from the catch model.

For $b = 1, \ldots, B$, compute a bootstrapped predicted total abundance $N_{tot}^{(b)}$ by adding up the predicted catches $y_k^{(b)}$ at each location $k = 1, \ldots, K$ as follows

1. Prediction of covariate values. The fitted covariate models for Secchi and conductivity are of the form $z_k \sim N(\hat{\mu}_k, \hat{\sigma}^2)$, where $\hat{\mu}_k = \mathbf{W}_k\hat{\boldsymbol{\beta}}$. $\mathbf{W}_k$ is a $1 \times J$ row vector of the design matrix with values corresponding to the intercept and soap film smooth basis of the latitude and longitude dimensions at location $k$, $\hat{\boldsymbol{\beta}}$ is a $J \times 1$ column vector of the estimated GAM parameter vector with $J \times J$ covariance matrix $\Sigma_{\hat{\boldsymbol{\beta}}}$, and $\hat{\sigma}^2$ is the estimated observation variance. For each covariate, obtain a $K \times 1$ column vector $\mathbf{z_b}$ of values as follows:

    (a) Simulate a $\boldsymbol{\beta}^{(b)} \sim N(\hat{\boldsymbol{\beta}}, \Sigma_{\hat{\boldsymbol{\beta}}})$.

    (b) Set $\mu_k^{(b)} = \mathbf{W}_k\boldsymbol{\beta}^{(b)}$

    (c) For $k = 1, \ldots, K$, simulate $z_k^{(b)} \sim N(\mu_k^{(b)}, \hat{\sigma}^2)$.

2. Construct the simulated covariate based design matrix $\mathbf{X}^{(b)}$ using the $\mathbf{z}^{(b)}$ values from step 1 for the Secchi and conductivity columns.

3. Prediction of catch given covariates. The catch model is of the form $y_k \sim NB(\hat{\lambda}_k, \hat{\theta})$, where $\log_e(\hat{\lambda}_k) = \mathbf{X}_k\hat{\boldsymbol{\beta}} + \log_e(Volume)$, $Vol_p$ is the volume sampled, $\hat{\theta}$ is the estimated dispersion parameter of the negative binomial distribution, $\mathbf{X}_k$ is the design matrix, $\hat{\boldsymbol{\beta}}$ is the estimated model parameter vector and $\Sigma_{\hat{\boldsymbol{\beta}}}$ is its covariance matrix. Given $\mathbf{X}^{(b)}$, $\log_e(\lambda_k, b)$ depends only on the value of a realization of $\boldsymbol{\beta}^{(b)} \sim N(\hat{\boldsymbol{\beta}}, \Sigma_{\hat{\boldsymbol{\beta}}})$. Viewed this way, $\log_e(\lambda_k^{(b)})$ are iid normal random variables with mean $\mathbf{X}_k\hat{\boldsymbol{\beta}} + \log_e(Volume)$ and variance

$$\tau_{\mu_k^{(b)}}^2 = \sum_{j=1}^{J} (x_{k,j}^{(b)})^2 \mathrm{Var}(\hat{\beta}_j) + \sum_{1 \le j < l \le J} 2x_{k,j}^{(b)} x_{k,l}^{(b)} \mathrm{Cov}(\hat{\beta}_j, \hat{\beta}_l)$$

7

where $x_{k,j}^{(b)} = \mathbf{X}_{k,j}^{(b)}$. To simulate catch values per $Vol_p$ of water,

(a) Simulate $\boldsymbol{\beta}^{(b)} \sim N(\hat{\boldsymbol{\beta}}, \Sigma_{\hat{\boldsymbol{\beta}}})$ and set $\mu_k^{(b)} = \mathbf{X}_k^{(b)}\boldsymbol{\beta}^{(b)} + \log_e(Vol_p)$.

(b) Compute the bias adjusted value $\mu_k^{(b,adj)} = \exp(\log_e(\mu_k^{(b)}) - \tau_{\mu_k^{(b)}}^2/2)$,

(c) For $k = 1, \ldots, K$, simulate $y_k^{(b)} \sim NB(\mu_k^{(b,adj)}, \hat{\theta})$

4. Compute subregion mean density, subregion abundance, and total abundance as

  (a) The mean catch density in subregion $h$ is

$$\bar{\delta}_h^{\ (b)} = \frac{\sum_{k \in h} y_k^{(b)}/K_h}{Vol_p}$$

where $K_h$ is the number of prediction locations in $h$.

(b) The total predicted number of fish per subregion is $\hat{N}_h^{(b)} = \bar{\delta}_h^{\ (b)}Vol_h$

  (c) The total predicted number of fish is

$$N_{tot}^{(b)} = \sum_{h=1}^{27} \hat{N}_h^{(b)}$$

Table S1: Abundance estimates and annual growth rates for February using design- and model-based approaches. Model based estimates are from Model 9 (Table 1 of the main text). Growth rates are year-over-year ratios.

| | Model-based | | Design-based | |
|---|---|---|---|---|
| Year | Abundance | Growth rate | Abundance | Growth rate |
| 2002 | 312488 | | 649028 | |
| 2003 | 268157 | 0.86 | 1479322 | 2.28 |
| 2004 | 269777 | 1.01 | 545347 | 0.37 |
| 2005 | 156633 | 0.58 | 291774 | 0.54 |
| 2006 | 91509 | 0.58 | 191509 | 0.66 |
| 2007 | 135563 | 1.48 | 100526 | 0.52 |
| 2008 | 43603 | 0.32 | 125398 | 1.25 |
| 2009 | 58877 | 1.35 | 142908 | 1.14 |
| 2010 | 55650 | 0.95 | 173163 | 1.21 |
| 2011 | 95682 | 1.72 | 86463 | 0.50 |
| 2012 | 271020 | 2.83 | 891304 | 10.31 |
| 2013 | 97707 | 0.36 | 74772 | 0.08 |
| 2014 | 56027 | 0.57 | 136596 | 1.83 |

Figure S5: Abundance estimates from model-based approaches using models M8, M9, and M12, and design-based. Grey shading shows the central 95% prediction interval for the M9 based predictions.

# References

Thompson, Steven K. 2002. Sampling, Second Edition. John Wiley & Sons, Inc. New York, NY

Wood, Simon W. 2006. Generalized Additive Models: An Introduction with R. Chapman & Hall, Boca Raton, FL

## Appendix

R code used to fit the models and compute predictions and prediction intervals. Complete
R code and input data available on request.

```
rm(list=ls())
library(MASS)
library(mgcv)
library(maptools)
library(proj4)
library(rgdal)
library(xtable)
library(rgeos)
library(car)
library(ncf)
library(geoR)
if(Sys.info()['sysname'] == 'Darwin'){
   library(parallel)
   }else{library(parallelsugar)}


data.root <- '~/smelt/gam-analyses/SKT-gam-analyses/Data/'

load(paste0(data.root,'SKT-2002-2014-gam-analysis-data-prep-v7.RData'))


# Time consuming model fits
load.M8.M12 <- TRUE
if(load(load.M8.M12)){
  load(file=paste0(data.root,'SKT-2002-2014-gam-analysis-soap-v9-M8.RData'))
  load(file=paste0(data.root,'SKT-2002-2014-gam-analysis-soap-v12-M8.RData'))
}

# M12
fit.M12 <- FALSE
if(fit.M12){
  m.regional.local.by.month.year.formula <- as.formula("smelt~offset(logVol)+
            fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb),by=fCohort.year.month)
                                              +Secchi.z+Cond.z+Tide")
  m.regional.local.by.month.year <- gam(m.regional.local.by.month.year.formula,
                                         family=nb(),
                                         knots=knots,data=ds,method="ML")
  save(list=ls(),
       file=paste0(data.root,'SKT-2002-2014-gam-analysis-soap-v9-M12.RData')
```

```
128   }
129
130   # M8
131   fit.M8 <- FALSE
132   if(fit.M8){
133     m.regional.by.month.year.formula <- as.formula("smelt~offset(logVol)+
134         fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb),by=fCohort.year.month)")
135     m.regional.by.month.year <- gam(m.regional.by.month.year.formula,family=nb(),
136                                     knots=knots,data=ds,method="ML")
137     save(list=ls(),
138         file=paste0(data.root,'SKT-2002-2014-gam-analysis-soap-v9-M8.RData')
139   }
140
141   # Remaining models
142   #VIF model
143   r.vif <- glm.nb(smelt~offset(logVol)+Secchi.z+Cond.z+Tide+Month+SubRegion,
144                   data=ds)
145   vif(r.vif)
146
147   r.vif <- glm.nb(smelt~offset(logVol)+Secchi.z+Cond.z+Tide+Month+Lon.z+Lat.z,
148                   data=ds)
149   vif(r.vif)
150
151   # 1) Global: no regional (smooth), no local
152   m.global <- glm.nb(smelt~offset(logVol)+fCohort.year*JD,data=ds)
153
154   # Global + one local to estimate the best a particular local covariate can do
155   m.global.plus.Secchi <- glm.nb(smelt~offset(logVol)+fCohort.year*JD+Secchi.z,data=ds)
156   m.global.plus.cond <- glm.nb(smelt~offset(logVol)+fCohort.year*JD+Cond.z,data=ds)
157   m.global.plus.tide <- glm.nb(smelt~offset(logVol)+fCohort.year*JD+Tide,data=ds)
158
159   # 2) Global x regional: no by in smooth
160   m.regional.formula <- as.formula("smelt~offset(logVol)+
161                               fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb))")
162   t1 <- Sys.time()
163   m.regional <- gam(m.regional.formula,family=nb(),
164                     knots=knots,data=ds,method="ML")
165   t2 <- Sys.time()
166   difftime(t2,t1)
167
168   # 3) Global x regional: by month
169   m.regional.by.month.formula <- as.formula("smelt~offset(logVol)+
170                     fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb),by=fmonth)")
171
172   t1 <- Sys.time()
```

```
173  m.regional.by.month <- gam(m.regional.by.month.formula,family=nb(),
174                              knots=knots,data=ds,method="ML")
175  t2 <- Sys.time()
176  difftime(t2,t1)
177
178
179  # 4) Global x regional: by year (cohort)
180  m.regional.by.year.formula <- as.formula("smelt~offset(logVol)+
181              fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb),by=fCohort.year)")
182
183  t1 <- Sys.time()
184  m.regional.by.year <- gam(m.regional.by.year.formula,family=nb(),
185                             knots=knots,data=ds,method="ML")
186  t2 <- Sys.time()
187  difftime(t2,t1)
188
189  # 6) Global x regional x local: No by
190  m.regional.local.formula <- as.formula("smelt~offset(logVol)+
191          fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb))+Secchi.z+Cond.z+Tide")
192
193  t1 <- Sys.time()
194  m.regional.local <- gam(m.regional.local.formula,family=nb(),
195                           knots=knots,data=ds,method="ML")
196  t2 <- Sys.time()
197  difftime(t2,t1)
198
199  # 7) Global x regional x local: by month
200  m.regional.local.by.month.formula <- as.formula("smelt~offset(logVol)+
201                      fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb),by=fmonth)
202                      +Secchi.z+Cond.z+Tide")
203
204  t1 <- Sys.time()
205  m.regional.local.by.month <- gam(m.regional.local.by.month.formula,
206                                    family=nb(),knots=knots,data=ds,method="ML")
207  t2 <- Sys.time()
208  difftime(t2,t1)
209
210  # 8) Global x regional x local: by year (cohort)
211  m.regional.local.by.year.formula <- as.formula("smelt~offset(logVol)+
212                  fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb),by=fCohort.year)
213                  +Secchi.z+Cond.z+Tide")
214
215  t1 <- Sys.time()
216  m.regional.local.by.year <- gam(m.regional.local.by.year.formula,family=nb(),
217                                   knots=knots,data=ds,method="ML")
```

12

```
218  t2 <- Sys.time()
219  difftime(t2,t1)
220
221
222  m <- m.regional.local
223
224  # Drop one local cov at a time to look at
225  # proportion deviance explained by adding this cov to a global X regional model
226  m.regional.local.minus.Secchi <- gam(smelt~offset(logVol)+
227                                       fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb))+
228                                       Cond.z+Tide,family=nb(),
229                                       sp=m$sp,knots=knots,data=ds,method="ML")
230
231  m.regional.local.minus.cond <- gam(smelt~offset(logVol)+
232                                     fCohort.year*JD+s(x,y,bs='so',xt=list(bnd=fsb))+
233                                     Secchi.z+Tide,family=nb(),
234                                     sp=m$sp,knots=knots,data=ds,method="ML")
235
236  m.regional.local.minus.tide <- gam(smelt~offset(logVol)+fCohort.year*JD+
237                                     s(x,y,bs='so',xt=list(bnd=fsb))+
238                                     Secchi.z+Cond.z,family=nb(),
239                                     sp=m$sp,knots=knots,data=ds,method="ML")
240
241  ##########--------- Make predictions of abundance in Feb ---------##########
242  # Make predictions at each grid location on Feb 15th
243  # for flood and ebb tides and bootstrap prediction intervals
244  ucym <- as.character(unique(ds$fCohort.year.month))
245  index.Feb <- which(as.numeric(unlist(lapply(as.character(ucym),
246                    FUN=function(x){y=strsplit(x,split="-")[[1]][3]})))==2)
247  ucym.Feb <- ucym[index.Feb]
248
249  # Make predictions on grid-
250  # Why doesn't crs(DSLCM.SubRegions) or CRS(DSLCM.SubRegions) work here?
251  # Why does crs work on a Windows PC? Or does it?
252  spatial.grid.predict <- SpatialPoints(grid.predict,
253                          proj4string=attributes(DSLCM.SubRegions)$proj4string)
254  grid.predict.with.subregions <- cbind(grid.predict,
255                                   over(spatial.grid.predict,DSLCM.SubRegions))
256
257  # Prediction volume
258  vol.p <- 10000
259
260  # Don't fix boundary at 0 for covariates
261  fsb.cov <- vector("list",1)
262  fsb.cov[[1]]$x <- region.boundary[,"x"]
```

```
263  fsb.cov[[1]]$y <- region.boundary[,"y"]
264
265  grid.cov.gam.func <- function(ucym,dat,cov){
266    # Get subregion averages of a given covariate for a dataset dt
267    dt <- subset(dat,fCohort.year.month==ucym)
268    names(dt)[which(names(dt)=="UTMX")] <- "x"
269    names(dt)[which(names(dt)=="UTMY")] <- "y"
270    cov.temp <- dt[,cov]
271    if(nrow(dt)>27){
272      cov.gam <- try(gam(cov.temp~s(x,y,bs="so",xt=list(bnd=fsb.cov)),
273                         knots=knots,data=dt,method="ML"))
274    }else{
275      cov.gam <- try(gam(cov.temp~s(x,y,k=25),data=dt,method="ML"))
276    }
277    return(cov.gam)
278  }
279
280  Secchi.z.models <- lapply(as.character(ucym),
281                            FUN=grid.cov.gam.func,dat=ds,cov='Secchi.z')
282  Cond.z.models <- lapply(as.character(ucym),
283                          FUN=grid.cov.gam.func,dat=ds,cov='Cond.z')
284
285  Secchi.z.gam.gof <- unlist(lapply(Secchi.z.models,
286                                    FUN=function(x){summary(x)$dev.expl}))
287  Cond.z.gam.gof <- unlist(lapply(Cond.z.models,
288                                  FUN=function(x){summary(x)$dev.expl}))
289
290  range(Secchi.z.gam.gof[index.Feb])
291  median(Secchi.z.gam.gof[index.Feb])
292  quantile(Secchi.z.gam.gof[index.Feb],probs=seq(.1,1,by=.1))
293  range(Cond.z.gam.gof[index.Feb])
294  median(Cond.z.gam.gof[index.Feb])
295  quantile(Cond.z.gam.gof[index.Feb],probs=seq(.1,1,by=.1))
296
297  grid.cov.gam.pred.func <- function(m,gcv.est=TRUE){
298    # Sample a prediction from a fitted GAM model m
299    # Returns a prediction at each location of the grid
300    # If gcv.est=F, prediction includes uncertainty in the
301    #    model coefficients and observation error
302    # Prediction does not include uncertainty in the smoothing paramter
303    data.new <- data.frame(x=grid.predict$x,y=grid.predict$y)
304    if(gcv.est){
305      y <- predict(m,newdata=data.new,type='response')
306    }else{
307      beta <- coef(m)
```

```
308     Vb <- m$Vc
309     Cv <- chol(Vb)
310     n.rep=1
311     nb <- length(beta)
312     br <- t(Cv) %*% matrix(rnorm(n.rep*nb),nb,n.rep) + beta
313     Xp <- predict(m,newdata=data.new,type="lpmatrix")
314     lp <- Xp %*% br
315     y <- rnorm(length(lp),mean=lp,sd=sqrt(m$sig2))
316   }
317   return(y)
318 }
319
320 cpue.newdata.grid.func <- function(fcym,tide.set,gcv.est){
321   index.temp <- which(ucym==fcym)
322   Secchi.z.temp <- grid.cov.gam.pred.func(m=Secchi.z.models[[index.temp]],
323                                           gcv.est=gcv.est)
324   Cond.z.temp <- grid.cov.gam.pred.func(m=Cond.z.models[[index.temp]],
325                                         gcv.est=gcv.est)
326
327   # Sets up a data frame of new data for making CPUE predictions
328   Month <- as.numeric(strsplit(as.character(fcym),split="-")[[1]][3])
329   if(Month==1){JD=15}
330   if(Month==2){JD=45}
331   if(Month==3){JD=74}
332   if(Month==4){JD=105}
333   if(Month==5){JD=135}
334
335   dn <- data.frame(
336     logVol=log(vol.p),
337     fCohort.year=factor(paste(
338       strsplit(as.character(fcym),split="-")[[1]][1:2],collapse="-"),
339       levels=levels(ds$fCohort.year)),
340     fCohort.year.month=factor(fcym,levels=levels(ds$fCohort.year.month)),
341     JD=JD,
342     fmonth=strsplit(fcym,split="-")[[1]][3],
343     Secchi.z=Secchi.z.temp,
344     Cond.z=Cond.z.temp,
345     Tide=tide.set,
346     x=grid.predict$x,
347     y=grid.predict$y
348   )
349   return(dn)
350 }
351
352 beta.param.vect.sample.from.gam <- function(m,b){
```

```r
353    # m a fitted GAM, returns a J x b column vector of
354    #  a samples of beta, J=length(beta)
355    beta <- coef(m)
356    Vb <- m$Vc
357    Cv <- chol(Vb)
358    n.rep <- b
359    nb <- length(beta)
360    br <- t(Cv) %*% matrix(rnorm(n.rep*nb),nb,n.rep) + beta
361    return(br)
362 }
363
364 var.sum.func <- function(a,x,Sigma){
365    # a, x- vectors of same length; Sigma- covariance matrix of x
366    # Let X=(a_1*x_1,...,a_n*x_n)
367    # Computes the variance of the of sum of the elements of X
368    # Var(sum(X))=sum_i a_i^2*Var(x_i)+2*sum_1<=i<j<=n a_i*a_j*Cov(x_i,x_j)
369
370    ai.aj <- combn(a,m=2,prod) # Get all a_i*a_j products for 1<=i<j<=n
371
372    # Get index of cov(x_i,x_j) values in same order as ai.aj vector
373    off.diag.index <- combn(1:ncol(Sigma),m=2) #All possible combinations of 1,...,n
374    off.diag.index <- cbind(off.diag.index[1,],off.diag.index[2,])
375
376    # Get the Cov(x_i,x_j) terms
377    off.diag.var.cov <- Sigma[off.diag.index]
378
379    # Compute variance of the sum
380    r <- sum(a^2*diag(Sigma))+sum(2*ai.aj*off.diag.var.cov)
381    return(r)
382 }
383
384 boot.pred.func <- function(m,fcym,tide.set,boot){
385    # ****Get covariance matrix of beta from m
386    # For bias correcting mu_k samples in the boot loop.
387    #   See **** Do this here for speed.
388    # Actual variance will depend on Xp.boot so need to wait till boot
389    #     loop to finish computing
390    off.diag.index <- combn(1:ncol(m$Vc),m=2) #All possible combinations of 1,..,n
391    off.diag.index <- cbind(off.diag.index[1,],off.diag.index[2,])
392    off.diag.var.cov <- m$Vc[off.diag.index]
393    v.m <- diag(m$Vc)
394
395    # Predictions at estimated parameters
396    data.new <- cpue.newdata.grid.func(fcym=fcym,tide.set=tide.set,gcv.est=T)
397    Xp <- predict(m,newdata=data.new,type="lpmatrix")
```

16

```
398    mu.pred.linear <- Xp %*% coef(m)+log(vol.p)
399    #GLM models don't bias correct when making prediction from log link models
400    mu.pred <- exp(mu.pred.linear)
401
402    mu.pred.mean.sr <- tapply(mu.pred,grid.predict.with.subregions$SubRegion,mean)
403
404    index.match <- match(wv$SubRegion,names(mu.pred.mean.sr))
405
406    tot.pop.size <- sum(mu.pred.mean.sr[index.match]*wv$twom,na.rm=T)/vol.p
407
408    if(boot==0){
409      return(list(
410        tot.pop.size=tot.pop.size,
411        mean.pop.boot=NA,
412        tot.pop.pred.boot.interval=NA
413      ))
414    }else{
415      boot.tot.pop.size <- rep(NA,boot)
416      theta.est <- m$family$getTheta(TRUE)
417      for(i in 1:boot){
418        # Step 1- simulate from covariate data models Xp_boot
419        boot.data.new <- cpue.newdata.grid.func(fcym=fcym,tide.set=tide.set,
420                                                 gcv.est=F)
421        Xp.boot <- predict(m,newdata=boot.data.new,type="lpmatrix")
422
423        # Sample a beta_b from N(hat(beta),Sigma_hat(beta))
424        beta.samp <- beta.param.vect.sample.from.gam(m=m,b=1)
425
426        # Make linear predictor using sample beta and sample covariates
427        boot.mu.pred.linear <- Xp.boot %*% beta.samp+log(vol.p)
428
429        # View log(tau_b)=boot.mu.pred.linear as a normally distributed variable
430        # tau_b ~ LN(logmean=Xp.boot %*% hat(beta)+log(vol.p),
431        #    varlog= Var(Xp.boot %*% hat(beta)+log(vol.p))=Var(Xp.boot %*% hat(beta))
432        # Then bias correct exp(log(tau_b))
433        # Bias correct assuming boot.mu ~ LN with mean=boot.mu.pred,
434        #    variance=sigma^2
435        # mu_i=1*beta_0+Xp1[i,1]*beta1+...+Xp[i,n]*beta_n
436        # mu=beta_0+x1*beta_1+x2*beta_2+...+xn*beta_n
437        # Var(mu)=sum_over_i x_i^2*var(beta_i)+
438        #    sum_over_i*sum_over_j x_i*x_j*Cov(beta_i,beta_j) covariance of sums formula
439        # Var(mu)=
440        #  sum_over_i x_i^2*var(beta_i)+2*sum_1<=i<j<=N x_i*x_j*cov(beta_i,beta_j)
441        # ****Have covariance matrix of beta
442        sig2 <- rep(NA,length(boot.mu.pred.linear))
```

17

```r
443        for(k in 1:length(boot.mu.pred.linear)){
444          ai.aj <- combn(Xp.boot[k,],m=2,prod) #All a_i*a_j products for 1<=i<j<=n
445          sig2[k] <- sum(Xp.boot[k,]^2*v.m)+sum(2*ai.aj*off.diag.var.cov)
446          #sig2[k] <- var.sum.func(a=Xp.boot[k,],x=coef(m),Sigma=m$Vc)
447        }
448        # End bias correct
449
450        boot.mu.pred <- exp(boot.mu.pred.linear-sig2/2)
451        boot.pred <- rnegbin(boot.mu.pred,theta=theta.est)
452        boot.dens.pred <- boot.pred/vol.p
453
454        boot.mean.dens.sr <- tapply(
455          boot.pred,grid.predict.with.subregions$SubRegion,mean)
456
457        boot.tot.pop.size[i] <- sum(
458          boot.mean.dens.sr[index.match]*wv$twom,na.rm=T)/vol.p
459      }
460      return(list(
461        tot.pop.size=tot.pop.size,
462        mean.pop.boot=mean(boot.tot.pop.size),
463        tot.pop.pred.boot.interval=quantile(
464          boot.tot.pop.size,probs=c(.025,.25,.5,.75,.975))
465      ))
466    }
467  }
468
469  # t1 <- Sys.time()
470  # e=boot.pred.func(m=m,fcym=ucym[2],tide.set="Flood",boot=2)
471  # t2 <- Sys.time()
472  # difftime(t2,t1)
473
474  # Check on understanding p1 should equal p1.alt
475  # t1 <- Sys.time()
476  # i=1
477  # fcym.temp <- as.character(ucym[i])
478  # data.new <- cpue.newdata.grid.func(fcym=fcym,tide.set=tide.set,gcv.est=T)
479  # p1 <- predict(m,newdata=data.new,type="response")
480  # Xp.alt <- predict(m,newdata=data.new,type="lpmatrix")
481  # p1.alt <- exp(Xp.alt %*% (coef(m))+log(vol.p))
482  # max(abs(p1-as.numeric(p1.alt)))
483  # t2 <- Sys.time()
484  # difftime(t2,t1)
485
486  # Point estimates and uncertainty using model 9
487  pop.estimate.posterior.sim.feb.func <- function(X,tide,boot){
```

```
488    r <- boot.pred.func(m=m,fcym=X,tide.set=tide,boot=boot)
489    return(r)
490  }
491
492  boot.set <- 4 #1000
493  cor.set <- 4
494
495  t1 <- Sys.time()
496  p.store.list.ebb.feb <- mclapply(X=ucym.Feb,
497                             FUN=pop.estimate.posterior.sim.feb.func,
498                             tide="Ebb",boot=boot.set,mc.cores=cor.set)
499  t2 <- Sys.time()
500  difftime(t2,t1)
501
502  boot.set <- 1000
503  t1 <- Sys.time()
504  p.store.list.flood.feb <- mclapply(X=ucym.Feb,
505                             FUN=pop.estimate.posterior.sim.feb.func,
506                             tide="Flood",boot=boot.set,mc.cores=cor.set)
507  t2 <- Sys.time()
508  difftime(t2,t1)
509
510  save(list=ls(),file=paste0(data.root,'SKT-2002-2014-gam-analysis-soap-v10.RData'))
511
512  # Point estimates using model 10
513  pop.estimate.posterior.sim.feb.func.M10 <- function(X,tide,boot){
514    r <- boot.pred.func(m=m.regional.local.by.month,fcym=X,tide.set=tide,boot=boot)
515    return(r)
516  }
517  p.store.list.flood.feb.M10 <- lapply(X=ucym.Feb,
518                             FUN=pop.estimate.posterior.sim.feb.func.M10,
519                             tide="Flood",boot=0)
520
521  # Point estimates using model 12
522  pop.estimate.posterior.sim.feb.func.M12 <- function(X,tide,boot){
523    r <- boot.pred.func(m=m.regional.local.by.month.year,fcym=X,tide.set=tide,
524                    boot=boot)
525    return(r)
526  }
527  p.store.list.flood.feb.M12 <- lapply(X=ucym.Feb,
528                             FUN=pop.estimate.posterior.sim.feb.func.M12,
529                             tide="Flood",boot=0)
530
531  pop.est.ebb.feb <- data.frame(
532    est=unlist(lapply(p.store.list.ebb.feb,FUN=function(x){x$tot.pop.size})),
```

```r
533    mean=unlist(lapply(p.store.list.ebb.feb,FUN=function(x){x$mean.pop.boot})),
534    median=unlist(lapply(p.store.list.ebb.feb,FUN=function(x){
535      x$tot.pop.pred.boot.interval['50%']})),
536    lower=unlist(lapply(p.store.list.ebb.feb,FUN=function(x){
537      x$tot.pop.pred.boot.interval['2.5%']})),
538    upper=unlist(lapply(p.store.list.ebb.feb,FUN=function(x){
539      x$tot.pop.pred.boot.interval['97.5%']})))
540  pop.est.flood.feb <- data.frame(
541    est=unlist(lapply(p.store.list.flood.feb,FUN=function(x){x$tot.pop.size})),
542    mean=unlist(lapply(p.store.list.flood.feb,FUN=function(x){x$mean.pop.boot})),
543    median=unlist(lapply(p.store.list.flood.feb,FUN=function(x){x$mean.pop.boot})),
544    lower=unlist(lapply(p.store.list.flood.feb,FUN=function(x){
545      x$tot.pop.pred.boot.interval['2.5%']})),
546    upper=unlist(lapply(p.store.list.flood.feb,FUN=function(x){
547      x$tot.pop.pred.boot.interval['97.5%']})))
548
549  pop.point.estimate.ebb.feb <- unlist(lapply(p.store.list.ebb.feb,
550                                      FUN=function(x){x$tot.pop.size}))
551  pop.point.estimate.flood.feb <- unlist(lapply(p.store.list.flood.feb,
552                                       FUN=function(x){x$tot.pop.size}))
553  pop.point.estimate.flood.feb.M10 <- unlist(lapply(p.store.list.flood.feb.M10,
554                                       FUN=function(x){x$tot.pop.size}))
555  pop.point.estimate.flood.feb.M12 <- unlist(lapply(p.store.list.flood.feb.M12,
556                                       FUN=function(x){x$tot.pop.size}))
557
558  plot(pop.point.estimate.flood.feb,
559       pop.point.estimate.flood.feb.M10,type='n',xlab='',ylab='')
560  abline(a=0,b=1)
561  title(xlab='Model 9 (single smooth for all months)',line=2.2)
562  title(ylab='Model 10 (month specific smooth)',line=2.2)
563  text(pop.point.estimate.flood.feb,pop.point.estimate.flood.feb.M10,
564       labels=sapply(uy,FUN=function(x){substr(x,start=3,stop=4)}))
565  cor(pop.point.estimate.flood.feb,pop.point.estimate.flood.feb.M10)
566  round(100*(pop.point.estimate.flood.feb.M10-pop.point.estimate.flood.feb)/
567        pop.point.estimate.flood.feb,2)
568
569  ab=data.frame(Year=uy,M9=pop.point.estimate.flood.feb,
570               M10=pop.point.estimate.flood.feb.M10,
571               M12=pop.point.estimate.flood.feb.M12)
572  print(ab,row.names=F)
573
574  ab=data.frame(
575    Coef=names(coef(m)[1:14]),
576    M9=coef(m)[1:14],
577    M10=coef(m.regional.local.by.month)[1:14],
```

```
578    M12=coef(m.regional.local.by.month.year)[1:14]
579  )
580  print(ab,row.names=F)
581
582  p.f <- function(x,p.t,col.set){
583    y.lim <- c(-10,3)
584    plot(x,type='n',ylim=y.lim,xaxt='n',xlab='',ylab='')
585    axis(side=1,at=x,labels=rownames(p.t),las=2,cex.axis=.85)
586
587    for(i in 1:nrow(p.t)){
588      points(x[i],p.t[i,'Estimate'],col=col.set,pch=20)
589      lines(rep(x[i],2),c(p.t[i,'Estimate']-p.t[i,'Std. Error'],
590                          p.t[i,'Estimate']+p.t[i,'Std. Error']),col=col.set)
591    }
592  }
593  p.f2 <- function(x,p.t,col.set){
594    for(i in 1:nrow(p.t)){
595      points(x[i],p.t[i,'Estimate'],col=col.set,pch=20)
596      lines(rep(x[i],2),c(p.t[i,'Estimate']-p.t[i,'Std. Error'],
597                          p.t[i,'Estimate']+p.t[i,'Std. Error']),col=col.set)
598    }
599  }
600
601  par(mar=c(10,3,2,1))
602  p.f(x=seq(1,31,by=1),p.t=summary(m.regional.local)$p.table,col.set='blue')
603  p.f2(x=seq(1.1,31.1,by=1),p.t=summary(m.regional.local.by.month)$p.table,
604       col.set='red')
605  p.f2(x=seq(1.2,31.2,by=1),p.t=summary(m.regional.local.by.month.year)$p.table,
606       col.set='green')
607  legend('bottomright',legend=c('M9','M10','M12'),
608         col=c('blue','red','green'),pch=20)
609
610  p.t.M9=summary(m)$p.table
611  p.t.M12=summary(m.regional.local.by.month.year)$p.table
612
613  delta0.M9 <- exp(c(p.t.M9[1,'Estimate'],
614                     p.t.M9[1,'Estimate']+p.t.M9[2:13,'Estimate']))
615  delta0.M12 <- exp(c(p.t.M12[1,'Estimate'],
616                     p.t.M12[1,'Estimate']+p.t.M12[2:13,'Estimate']))
617
618  print(data.frame(cohort=cohort,M9.delta0=delta0.M9,
619                   M12.delta0=delta0.M12,ratio=delta0.M12/delta0.M9),row.names=F)
620
621  data.frame(cohort=cohort,M9.delta0=delta0.M9,
622             M12.delta0=delta0.M12,ratio=delta0.M12/delta0.M9)
```

```r
623  s.t.M9=summary(m)$s.table
624  s.t.M12=summary(m.regional.local.by.month.year)$s.table
625
626  a=names(coef(m.regional.local.by.month.year))
627
628
629  # Summaries across models
630  # Theta estimates
631  theta.est <- c(m.global$theta,
632  m.global.plus.Secchi$theta,
633  m.global.plus.cond$theta,
634  m.global.plus.tide$theta,
635  unlist(lapply(list(m.regional,
636     m.regional.by.month,
637     m.regional.by.year,
638     m.regional.by.month.year,
639                 m.regional.local,
640                 m.regional.local.by.month,
641                 m.regional.local.by.year,
642          m.regional.local.by.month.year,
643                 m.regional.local.minus.Secchi,
644                 m.regional.local.minus.cond,
645                 m.regional.local.minus.tide),FUN=function(x){
646                   return(x$family$getTheta(TRUE))}))
647  )
648  theta.est
649  range(theta.est)
650  m$family$getTheta(TRUE)
651
652  # Model comparison
653  prop.dev.func <- function(a){
654  return((a$null.deviance-a$deviance)/a$null.deviance)
655  }
656
657  AIC.set <- AIC(m.global,
658     m.global.plus.Secchi,
659     m.global.plus.cond,
660     m.global.plus.tide,
661     m.regional,
662     m.regional.by.month,
663     m.regional.by.year,
664     m.regional.by.month.year,
665                 m.regional.local,
666                 m.regional.local.by.month,
667                 m.regional.local.by.year,
```

```
668          m.regional.local.by.month.year,
669              m.regional.local.minus.Secchi,
670              m.regional.local.minus.cond,
671              m.regional.local.minus.tide
672              )
673
674 index.temp <- which.min(AIC.set$AIC)
675 delta.AIC.set <- AIC.set$AIC-AIC.set$AIC[index.temp]
676
677 ML.score.set <- c(NA,NA,NA,NA,unlist(lapply(list(
678    m.regional,
679    m.regional.by.month,
680    m.regional.by.year,
681    m.regional.by.month.year,
682              m.regional.local,
683              m.regional.local.by.month,
684              m.regional.local.by.year,
685          m.regional.local.by.month.year,
686              m.regional.local.minus.Secchi,
687              m.regional.local.minus.cond,
688              m.regional.local.minus.tide),FUN=function(x){x$gcv.ubre})))
689
690 Percent.dev.exl.set <- unlist(lapply(list(
691    m.global,
692    m.global.plus.Secchi,
693    m.global.plus.cond,
694    m.global.plus.tide,
695    m.regional,
696    m.regional.by.month,
697    m.regional.by.year,
698    m.regional.by.month.year,
699              m.regional.local,
700              m.regional.local.by.month,
701              m.regional.local.by.year,
702          m.regional.local.by.month.year,
703              m.regional.local.minus.Secchi,
704              m.regional.local.minus.cond,
705              m.regional.local.minus.tide),FUN=function(x){prop.dev.func(x)}))
706
707 dev.expl.table <- data.frame(
708 Model=c('Global',
709 'Global + Secchi',
710 'Global + cond',
711 'Global + tide',
712 'Global + regional (single)',
```

```
713  'Global + regional (month)',
714  'Global + regional (year)',
715  'Global + regional (month year)',
716          'Global + regional (single) + local',
717          'Global + regional (month) + local',
718          'Global + regional (year) + local',
719      'Global + regional (year) + local',
720          'Global + regional (single) + local - Secchi',
721          'Global + regional (single) + local - cond',
722          'Global + regional (single) + local - Tide'),
723  Df=AIC.set$df,
724  AIC=AIC.set$AIC,
725  delta.AIC=delta.AIC.set,
726  SSC=ML.score.set,
727  'Dev. exp.'=100*Percent.dev.exl.set,
728  'Theta'=theta.est
729  )
730  dev.expl.table <- data.frame(M=1:nrow(AIC.set),Model=dev.expl.table$Model,round(dev.expl
731  dev.expl.table
732  write.csv(dev.expl.table,
733          file='~/smelt/gam-analyses/SKT-gam-analyses/dev-explained-table-v10.csv',
734          row.names=F)
735  print.xtable(xtable(dev.expl.table,digits=1),include.rownames=F)
```