

# St. Pierre bank assessment model

Jonathan Babyn

March 25, 2019

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Model Design</b>	<b>2</b>
2.1	Process Equations . . . . .	2
2.2	Observation Equations . . . . .	3
<b>3</b>	<b>Data Inputs</b>	<b>5</b>
3.1	Smoothing Maturities and Weights . . . . .	6
<b>4</b>	<b>Model Analysis</b>	<b>7</b>
	<b>Appendices</b>	<b>7</b>
<b>A</b>	<b>Numerically stable logged normal cumulative distribution function derivatives in TMB</b>	<b>7</b>
A.1	Extending stability to censored bounds . . . . .	8
<b>B</b>	<b>Functional Principal Component Analysis</b>	<b>9</b>
	<b>Appendices</b>	<b>10</b>
<b>C</b>	<b>C++ Code used</b>	<b>10</b>
C.1	fit.cpp (Model Main) . . . . .	10
C.2	pnorm4.hpp . . . . .	18
C.3	Fcorr.hpp . . . . .	21
C.4	MVMIX.hpp . . . . .	23
C.5	C++ code used outside model . . . . .	25
C.6	R Code used . . . . .	26
C.6.1	spam.R . . . . .	26
C.6.2	NCplots.R . . . . .	32
C.6.3	fpcaData.R . . . . .	39
C.6.4	mdata3Ps.R . . . . .	43

# 1 Introduction

The cod stock in Northwest Atlantic Fisheries Organization (NAFO) subdivision 3Ps is the last commercial cod fishery open in Newfoundland. NAFO subdivision 3Ps is located off the southern coast of Newfoundland containing the French Overseas Collectivity of Saint Pierre & Miquelon is located partially on St. Pierre bank. The cod stock there is jointly managed by the French and Canadian governments. The stock is currently considered to be in the cautious zone by DFO as of January 2017[[rideout2017assessing](#)].

Currently the stock assessment model for 3Ps is a survey based model referred to as SURBA written by Noel Cadigan that has been in use since 2010[[cadigan2010trends](#)]. Since SURBA is a survey based model it only incorporates information from surveys such as Fisheries and Oceans Canada (DFO)'s multispecies Research Vessel (RV) bottom trawl survey that has occurred almost every year in 3Ps since 1972[[doubleday1981manual](#)]. The current 3Ps stock assessment model also relies on data solely from DFO's RV survey from 1983 onwards and not from the French survey that was performed from 1978 to 1992, the industry backed survey by the Groundfish Enterprise Allocation Council (GEAC) that occurred in the area from 1997 to 2007, or the still ongoing sentinel survey that started in 1995. Some of these surveys have been involved in past assessment models but have not been used recently[[bratney2005assessment](#)].

## 2 Model Design

### 2.1 Process Equations

The St. Pierre bank Assessment Model (SPAM) is a state-space stock assessment model implemented in Template Model Builder (TMB). As a State Space Model (SSM) the model can be separated into process equations and observation equations. SPAM follows the same general process equations design as Nielsen and Berg[[Nielsen2014Estimation-of-t](#)], with the log abundance at age  $a$  in a given year  $y$ ,  $\log(N_{y,a})$  as

$$\log(N_{y,R}) = \log(\text{SR}) + \eta_{y,R}, \quad \eta_{y,R} \sim \mathcal{N}(0, \sigma_R^2) \quad R < a, \quad (1)$$

$$\log(N_{y,a}) = \log(N_{y-1,a-1} \exp(-Z_{y-1,a-1})) + \eta_{y,a}, \quad \eta_{y,a} \sim \mathcal{N}(0, \sigma_S^2) \quad R < a < A \quad \text{and} \quad (2)$$

$$\log(N_{y,A}) = \log(N_{y-1,a-1} \exp(-Z_{y-1,a-1}) + N_{y-1,A} \exp(-Z_{A,y-1})) + \eta_{y,A} \quad \eta_{y,A} \sim \mathcal{N}(0, \sigma_S^2) \quad (3)$$

with  $Z_{y,a} = F_{y,a} + M_{y,a}$ , where  $F_{y,a}$  is mortality from fishing and natural mortality  $M_{y,a}$ . Natural mortality was taken to be a fixed constant.

Equation 1 is the process of recruitment of cod to the fishery,  $R$  is the first age included in the model and SR is some form of stock recruitment. SPAM implements four different options for stock recruitment,

$$\text{Random Walk: } \text{SR} = N_{y,R} = N_{y-1,R} + \eta_{y,R} \quad (4)$$

$$\text{Ricker: } \text{SR} = \alpha \text{SSB} \exp(-\beta \text{SSB}) \quad (5)$$

$$\text{Beverton-Holt: } \text{SR} = \frac{\alpha \text{SSB}}{1 + \beta(\text{SSB})} \quad \text{and} \quad (6)$$

$$\text{Smooth HS: } \text{SR} = \alpha \left( \text{SSB} + \sqrt{\delta^2 + \frac{\gamma^2}{4}} - \sqrt{(\text{SSB} - \delta)^2 + \frac{\gamma^2}{4}} \right) \quad (7)$$

where SSB is Spawning Stock Biomass (SSB) and  $\gamma^2$  is a fixed parameter that controls the smoothness of the breakpoint and was set equal to 0.1. Equations 2 & 3 are the process of survival, and

Name	Off-diagonal elements
Independent	$\Sigma_{a,\tilde{a}} = 0$
Parallel	$\Sigma_{a,\tilde{a}} = \sigma_a \sigma_{\tilde{a}}$
Compound	$\Sigma_{a,\tilde{a}} = \rho \sigma_a \sigma_{\tilde{a}}$
AR(1)	$\Sigma_{a,\tilde{a}} = \rho^{ a-\tilde{a} } \sigma_a \sigma_{\tilde{a}}$
Custom	$\Sigma_{a,\tilde{a}} = \rho_a \rho_{\tilde{a}} \sigma_a \sigma_{\tilde{a}}$

Table 1: Covariance matrix options available in SPAM for ages  $a$ ,  $a \neq \tilde{a}$ , independent is uncorrelated between ages and the process will develop independently in time, parallel is like  $\rho = 1$  and is the same as the assumption of constant selectivity, compound has all the age groups develop correlated in time the same, AR(1) has all the processes be correlated as function of the distance between ages, so closer ages are more closely related. Custom allows having ages have their own correlation parameter.

Equation 3 is an extension for a plus group of cod aged  $A$  or older. The entire  $\mathbf{N}$  abundance matrix is treated as random parameters to be integrated out of the model using the Laplace Approximation, this has the benefit of making adding internal stock recruitment models very simple as the calculation of SSB can be done at any point in the model, as the values of  $\mathbf{N}$  will be found during optimization. It also greatly reduces model run-time and the number of fixed effect parameters to estimate. Modelling survival as process error also has the benefit of allowing for immigration of fish into the stock of interest, which is something that may be relevant to 3Ps cod as stock mixing has been a concern[[methot2005spatio](#)][[rideout2017assessing](#)].

Fishing mortality was also implemented the same way as in Nielsen & Berg[[Nielsen2014Estimation-of-t](#)]. For a given year of fishing mortality  $\mathbf{F}_y = (F_{y,1}, \dots, F_{y,A})'$  is a Multivariate Normal (MVN) random walk,

$$\log(\mathbf{F}_y) = \log(\mathbf{F}_{y-1}) + \xi_y, \quad \xi_y \sim MVN(\mathbf{0}, \Sigma_{\mathbf{F}}) \quad (8)$$

, with  $\Sigma_{\mathbf{F}}$  being the covariance matrix. This implementation allows for time varying selectivity, with selectivity being defined as  $S_{y,a} = \frac{F_{y,a}}{\sum_a F_{a,y}}$ [[Nielsen2014Estimation-of-t](#)]. Since the 3Ps cod fishery has seen a shift in gear types over time, most notably a move away from otter trawls and cod traps to gillnet, time varying selectivity can help deal with this. As in Nielsen and Berg[[Nielsen2014Estimation-of-t](#)], the same four covariance options were implemented. A fifth covariance option was also added which allows for a different correlation  $\rho_a$  for each age. This was mainly done with the intention of making fishing mortality in age 2 cod uncorrelated with older fish, e.g.  $\rho_2 = 0$ . Table 1 illustrates the different covariance matrix options available in SPAM.

## 2.2 Observation Equations

SPAM incorporates data from multiple different research vessel surveys, catch at age proportions and reported landings. For the survey indices  $I_{y,a,s}$ , which are reported as Mean Numbers Per Tow (MNPT), the value predicted by the model is

$$\log(\hat{I}_{y,a,s}) = \log(q_{a,s}) + \log(N_{y,a}) - t_{s,y} Z_{y,a} \quad (9)$$

where  $q_{a,s}$  is the catch-ability parameter for survey  $s$  for age  $a$ , and  $t_{s,y}$  is the median date of the survey during the year in decimal form. As in the Northern Cod Assessment Model (NCAM) currently used for assessment for northern cod, SPAM sets a survey detection limit of 0.0005, with the belief that the cod are there, just in not large enough numbers to be detected by the

survey[Cadigan2016A-state-space-s]. This says that the data is left censored.  $I_{y,a,s} > 0.0005$  are added to the log-likelihood as

$$\sum_a \sum_y \log \left( \varphi \left( \frac{\log(I_{y,a,s}) - \log(\hat{I}_{y,a,s})}{\sigma_{s,I}} \right) \right) \quad (10)$$

where  $\varphi$  is the PDF of the standard normal distribution and  $\sigma_{s,I}$  is the standard deviation of the survey index for survey  $s$ . Indices that fall below the detection limit are added to the likelihood as

$$\sum_a \sum_y \log \left( \Phi \left( \frac{\log(I_{y,a,s}) - \log(\hat{I}_{y,a,s})}{\sigma_{s,I}} \right) \right) \quad (11)$$

with  $\Phi$  being the Cumulative Distribution Function (CDF) of the standard normal distribution.

SPAM calculates the expected catch  $C$  in the model using the Baranov catch equation,

$$\hat{C}_{y,a} = N_{y,a}(1 - \exp(-Z_{y,a})) \frac{F_{y,a}}{Z_{y,a}}. \quad (12)$$

It does not fit directly to  $\hat{C}$  but instead to the expected Continuation Ratio Logit (CRL) of the catch proportions at age for a given year,  $\hat{\mathbf{X}}_y$ . CRL are defined as  $P(Y = a|a \geq a) = \text{logit}(\pi_a)$  where  $\pi_a = \frac{p_a}{p_a + \dots + p_A}$  and  $p_a$  is proportion of catch at age  $a$ . Using CRLs allows for fitting on catch proportions but does not require knowing or assuming the sampling sizes used to determine the catch proportions. It does require having non-zero catch proportions and small values were added to the observed zero catch proportions, and these were used to create the observed CRL,  $\mathbf{X}_y$ . The observed CRL are assumed be distributed as

$$\mathbf{X}_y \sim MVN(\hat{\mathbf{X}}_y, \Sigma_{CRL}) \quad (13)$$

where  $\Sigma_{CRL}$  can be any one of the same covariance matrix options presented in Table 1 allowing for a correlation structure between ages.

Since SPAM fits catch proportions instead of catch at age, a third element is needed to fix the the magnitudes of abundance. The reported landings that were provided were used to do this. This was done in a similar fashion to how it is outlined in the NCAM paper, except instead of fitting to expected total catch for the year versus the observed total catch, this was instead done on the expected landings versus the observed landings. This was done since earlier catch at age calculations have been speculated to not include fish older than age 14 rather than their just be zero catch at those ages, whereas reported landings would and thus total observed catch would be under-reporting. The landings are added to the log-likelihood in the model as

$$\sum_y \log \left\{ \Phi \left( \frac{\log(L_y) + \log(UB) - \log(\hat{L}_y)}{\sigma_L} \right) - \Phi \left( \frac{\log(L_y) + \log(LB) - \log(\hat{L}_y)}{\sigma_L} \right) \right\} \quad (14)$$

where  $L_y$  and  $\hat{L}_y$  are the observed and predicted landings respectively,  $UB$  and  $LB$  are upper and lower bound multipliers on what the reported landings actually were. This is similar to Equation 7 of the NCAM paper. There it was noted that there can be issues optimizing the likelihood with this sort of formulation and there a mixture distribution and judicious choices of starting parameters were used to overcome them[Cadigan2016A-state-space-s]. I instead opted to fix the problem directly by writing atomic functions for TMB that manually defines the derivatives for Equations 11 & 14 in a more numerically stable manner. This process is described in Appendix A. This

Fixed Parameters	Description
$\sigma_R$	Std. Dev. of Recruitment Process
$\sigma_S$	Std. Dev. of Survival Process
$\sigma_{s,I}$	Vector of Std. Dev. of survey indices
$\sigma_F$	Vector of Std. Dev. for fishing mortality
$\rho_F(\rho_F)$	Correlation parameter for $\Sigma_F$ (Vector if custom type, 0 if independent)
$\sigma_{CRL}$	Vector of Std. Dev. for CRLs.
$\rho_{CRL}(\rho_{CRL})$	Correlation parameter for $\Sigma_{CRL}$ (Vector if custom type, 0 if independent)
$\mathbf{q}$	Vector of catchability parameters for the survey indices
$\alpha$	Optional parameter used in the three stock recruitment functions
$\beta$	Optional parameter used in the Ricker and Beverton-Holt stock recruitment functions
$\delta$	Optional parameter used in the smooth hockey stick recruitment function
Random Parameters	
$\mathbf{N}$	Matrix of Abundance
$\mathbf{F}$	Matrix of Fishing Mortality

Table 2: Mandatory and optional fixed effect parameters estimated in SPAM followed by random parameters estimated.

allows for SPAM to be run with a wide variety of starting values and without the need for upper or lower bounds on parameters or the need to run the model twice to find reasonable starting values. For the same reason outlined in Cadigan (2016)[**Cadigan2016A-state-space-s**],  $\sigma_L$  is fixed to be  $\log(0.02)$  to prevent  $\hat{L}_y$  from escaping the bounds too often or exceeding them by too much.

The parameters estimated, including optional parameters by SPAM are noted in Table 2.

### 3 Data Inputs

Catch, surveys, landings, weights and maturities were used from 1959 to 2016 if available. The assessment here is run on data ranging from ages 2 to 16 with 14 and up comprising a plus group. For the plus group I used the sum for the catch and survey indices for ages 14 to 16 and the mean for the weights and maturities.

The DFO supplied RV surveys started in 1983 and continue to run to this day. In 1997 the survey was extended to add more inshore strata. I decided to treat the DFO survey from 1983-1996 with offshore indices only as one survey, and the 1997 onwards inshore plus offshore strata as another. In addition to the DFO Offshore and Inshore and Offshore survey indices that were provided, I also took the RV survey indices that were the result of the Groundfish Enterprise Allocation Council (GEAC) RV survey performed in 3Ps. GEAC performed their survey in 3Ps from 1997 to 2007, with 2006 excluded and their data was given to DFO to be analyzed in the same manner as DFO's RV survey. These survey indices were taken from the DFO research documents in [**mcclintock2005year**] & [**mcclintock2011fall**] and made available for everyone else in the class to use. The GEAC indices were included in past 3Ps cod assessments, but have not been included recently, in addition the 1997 survey indices were not used in those assessments due to a much smaller area being surveyed than other years[**bratney2005assessment**]. Due to this I also elected not to use the 1997 GEAC survey indices and I also chose not to use the 2007 GEAC survey indices since they expanded the survey extent and used an entirely different set of equipment than previous years. Juliette Champagnat provided the survey indices from the French RV survey that was performed in 3Ps from 1978 to 1992.

Due to the approach I opted to take with the landings, upper and lower bounds on the landings are required. I consulted with Danny Ings, head of 3Ps cod stock assessment at DFO, he

Period	Danny Ings' Description	Tight	Medium	Wide
1959-1993	Somewhat wide	(0.75,1.25)	(0.5,1.5)	(0.35,1.75)
1994-1996	Narrow	(0.95,1.05)	(0.90,1.10)	(0.80,1.20)
1997-1998(or 1997-1999)	Very Wide	(0.5,1.5)	(0.35,1.75)	(0.25,2.0)
1999-present(or 2000-present)	Narrower than 1959-1993 but wider than 1994-96	(0.85,1.15)	(0.75,1.25)	(0.5,1.5)

Table 3: The periods for landing bounds provided by Danny Ings, along with the description of how they relate and three sets of different upper and lower bound multipliers to be tested on SPAM.

provided a set of periods with a general idea of how they relate to each other in width, but not exact bounds. These periods, along with the selection of landings bounds that were tried can be seen in Table 3.

As previously mentioned, since the raw catch proportions used to determine catch at age are not available, nor their sample sizes, the catch at age was converted to be catch proportions at age then to the catch CRLs used in the model. Ideally the raw catch proportions would be provided and their sample sizes known which would allow for other possibly better options of representing compositional data.

### 3.1 Smoothing Maturities and Weights

Two important pieces for the model are the weights used and maturity ogives. The maturity ogives provided are the output of a logistic regression run and updated each year from fish that have had their otoliths aged, lengths measured and maturity status determined. This results in very noisy maturity data set for each age group. For both the commercial and stock weights, output from a cohort model going from 1983 to 2016 were provided, along with commercial and stock weights from 1959 to 2014 and 1959 to 2017 respectively provided by DFO. The weights provided by DFO are noisy, and values prior to 1977 are identical and likely copy and pasted from some mean weights derived at some point. Neither of these situations are ideal. The weights either need to be extended thirty back into the past or improved. While the DFO provided weights unfortunately do not have real observations prior to 1977, the values provided at least seem consistent when compared with the output of the cohort model weights. I chose to try and deal with the noise in the DFO provided weights rather than extend the other weights back as the DFO weights are based on actual commercial sampling efforts rather than being based off research vessel survey data.

Both the maturities and weights were smoothed using Functional Principal Component Analysis (FPCA). FPCA is a non-parametric technique that has been used to perform noise reduction on curve like data[ramsay2007applied]. FPCA has at least been at least applied to fisheries data in the context of evaluating spatio-temporal patterns in the data[embling'2012] but it has not to the author's knowledge been applied to fisheries weight or growth data. The noise reduction is done by breaking the data down into it's mean function  $\mu(t)$ , eigenfunctions  $\phi_k(t)$  and Functional Principal Components (FPC),  $\xi_{ki}$  and keeping only the number of FPCs that explain most of the variation in the data. Further details on FPCA and this process can be seen in Appendix B.

For the maturities, the logit was taken and the FPCA performed in logit space to ensure that maturity proportions would continue to lie between 0 and 1 after smoothing. The FPCA was performed with each age group being considered an observation measured at each year. Smoothing was allowed on the mean function and a user selected bandwidth of 14 was chosen which removed many of the peaks and valleys. The diagnostic plots of the FPCA performed on the maturities showcasing the mean function, first three eigenfunctions, fraction of variance explained by each FPC are seen in Figure ?? . Over 90% of the variance is explained by the first FPC, because of this only the first FPC was used in creating the smoothed data. Plots comparing the smoothed and

raw versions of the maturities by age and year are in Figure ??.

For the weights the years were treated as observations and the age groups as measurement points. Here the mean function was automatically smoothed using Generalized Cross Validation (GCV). The diagnostic plots for the FPCA on the commercial weights are in Figure ?. As the commercial weights represent mid-year weights, the points were placed on a grid at half year intervals, e.g. the weights for fish that are aged 3.5 years. Looking at the smoothed commercial weights I noticed that the main bulk of the data appeared to be broken into two groups and could be mostly split into data pre and post 1994. This trend is also noticable in the original data but not as clearly. There was a shift in gear types towards gillnets in 3Ps in 1994, however that should result in the opposite trend with commercial weights increasing if this was the cause. Low sampling numbers would most likely not explain this either, as a 20+ year sustained downward bias would be unlikely. There is the possibility that there has been a change in how the aging has been done either with pre-1994 cod otolith readers underestimating ages or post-1994 readers overestimating ages. This change also occurs during the period when the stock collapsed and it is possible that this is somehow related as well. Danny Ings told me DFO is aware of the issue but is not sure of the exact cause either when I consulted him about it.

Danny Ings also informed me that the provided DFO stock weights are just model back calculations of the commercial weights. I decided to move the smoothed commercial weights back half of a year in time and use those as the stock weights for the model. This is effectively just shifting them down the curve.

## 4 Model Analysis

There are too many possible options in SPAM to evaluate all of them directly. To narrow down a set of candidate model options, I ran SPAM once for every possible option of covariance matrix for the  $\mathbf{F}_y$  random walk and the catch CRLs along with the four possible recruitment options on a single specified mapping of vector parameters shown in table ?. The Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC), log-likelihood, time elapsed and convergence status were recorded. 66 of the 100 model runs resulted in relative convergence while the remaining 44 resulted in false or singular convergence and were removed from contention. The top performing models were those that used the custom correlation structure for both the  $\mathbf{F}_y$  random walk and the catch CRLs, note that  $\rho_2$  was explicitly fixed to be equal to 0 in both cases.

# Appendices

## A Numerically stable logged normal cumulative distribution function derivatives in TMB

TMB uses Automatic Differentiation (AD) to generate the first and second derivatives which are used to help with both performing MLE for parameters and integrating random effects out of the likelihood with the Laplace Approximation (LA). TMB accomplishes this with a mixture of forward and reverse AD accumulation modes. Further details on AD can be seen in Fournier et al[**Fournier’2012**]. TMB automatically tries to determine the derivatives of the user’s template. However, it can not automatically determine the most numerically stable form of the derivative which can often be necessary when working with extremely small values that can commonly occur when working with probabilities. TMB also includes the ability to manually define the forward and

reverse modes for a given function to overcome this limitation using atomic functions. By default TMB includes a built in atomic version of the normal CDF, `pnorm` based on the version in R's C language math library. This built-in version does not support the ability to return a logged value of the probability and shares the same limitation of returning 0 or 1 for negative or positive values of standard normal  $Z$  scores when  $|Z| > \sim 38$  due to floating point limitations. This clearly means that for any  $Z < -38$ , TMB will return a NaN in the gradient when trying to perform  $\log(\text{pnorm}(Z))$  due to  $\log(0)$  being negative infinity. In practice however this is actually much worse since the derivative of  $\log(\text{pnorm}(Z))$  is

$$\frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi} \left( \frac{\text{erf}\left(\frac{z}{\sqrt{2}}\right)}{2} + \frac{1}{2} \right)} \quad (15)$$

and easily results in the division of two extremely small numbers. With floating point limitations this quickly becomes undefined if not handled correctly.

I created an atomic function to specifically deal with the problem of numerical stability when trying to do  $\log(\text{pnorm}(Z))$  and similar calculations in the objective function of SPAM such as those done when using a censored likelihood for the detection limit for fitting the observation equations for the survey indices. This requires specifying more numerically stable versions for both the forward and reverse AD modes. For the forward mode this is just the `double C++` function giving the log of the normal CDF and all that requires is calling the version of `pnorm` in R's C math library with the `give.log` flag turned on. For the reverse mode just using Equation 15 is not numerically stable due to the division. It can easily be seen that Equation 15 is equivalent to

$$\exp \left( \log \left( \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} \right) - \log \left( \frac{\text{erf}\left(\frac{z}{\sqrt{2}}\right)}{2} + \frac{1}{2} \right) \right). \quad (16)$$

Working in log space is much less likely to result in under or overflow of floating point values when performing division of small numbers. Using numerically stable versions of the logged normal PDF & CDFs and the form in Equation 16 results in a numerically stable reverse mode derivative for  $\log(\text{pnorm}(Z))$ . The atomic function was wrapped in the function `pnorm4` and made available for others to use in the C++ header file `pnorm4.hpp`. Using `pnorm4(Z)` in place of  $\log(\text{pnorm}(Z))$  allows for running one-sided censor likelihood bounds without the need for using two runs of optimization.

It's accuracy was checked by comparing the results of the gradient and hessian generated by TMB for `pnorm4` against the numerical first and second derivatives of `pnorm4` done by the R software package `numDeriv` along with the brute force first and second derivatives of the normal  $\log(\text{CDF})$  calculated by the open-source computer algebra system `Maxima` with 5000 digits of floating point precision.

## A.1 Extending stability to censored bounds

Censored likelihood bounds pose a similar problem to the above. When using the built-in functions to add normally distributed censored log-likelihood bounds this can be done like  $\log(\text{pnorm}(Z_U)) - \log(\text{pnorm}(Z_L))$  where  $Z_U$  &  $Z_L$  are  $Z$  scores of the upper and lower bounds respectively. Just replacing the calls to `pnorm` with the more stable `pnorm4` and using the brute force approach or `logspace_sub` to perform the subtraction is not stable either. This is because for  $Z > 38$  `pnorm` will return 1, or 0 for the logged version which will then result in trying to take the log of 0, again leading to NaN in the gradient.



The normal CDF of the upper and lower bounds can be thought of as  $\Phi(ZU) = 1 - \epsilon_u$  and  $\Phi(ZL) = 1 - \epsilon_l$  where the  $0 < \epsilon < 1$  and the censored log-likelihood bound is

$$\log(\Phi(ZU) - \Phi(ZL)) = \log(1 - \epsilon_u - (1 - \epsilon_l)) = \log(\epsilon_l - \epsilon_u). \quad (17)$$

It can also be seen that

$$\log(\Phi(-ZL) - \Phi(-ZU)) = \log(\epsilon_l - \epsilon_u) \quad \text{since} \quad \Phi(-Z) = 1 - \Phi(Z) = 1 - (1 - \epsilon) = \epsilon. \quad (18)$$

Since the problem of NaNs in the gradient is occurs because the relatively large value of 1 overwhelms the incredibly small values of  $\epsilon$  for large positive values of  $Z$  again due to floating point limitations. This was solved by using equation 18 for positive values of  $Z$  instead but with probabilities returned in logged form and the subtraction done with `logspace_sub`. This prevents the small  $\epsilon$ s from disappearing in the subtraction.

Then a stable reverse mode similar to the one for `pnorm4` was made, and this was wrapped in the function `censored_bounds`, checked for accuracy the same way and again made available for everyone to use in `pnorm4.hpp`.

## B Functional Principal Component Analysis

Functional Data Analysis (FDA) is a branch of statistics that deals with data that takes values in an infinite dimensional or functional space. Growth and maturity curves are one-dimensional examples of functional data. Several classical statistical techniques have been extended to work with functional data, Principal Component Analysis (PCA) being among them. In Principal Component Analysis (PCA) the goal is to perform dimension reduction while maximizing the variation explained by each dimension, Functional Principal Component Analysis (FPCA) extends this to functional data. Let  $X$  be a random function defined on the function grid  $[0, T]$  with unknown smooth mean function

$$E[X(t)] = \mu(t) \quad t \in [0, T], \quad (19)$$

and covariance function

$$\text{Cov}(X(s), X(t)) = G(s, t) \quad s, t \in [0, T]. \quad (20)$$

$G(s, t)$  can be written in it's orthogonal expansion as

$$G(s, t) = \sum_{k=1}^{\infty} \lambda_k \phi_k(s) \phi_k(t) \quad (21)$$

where  $\lambda_k$  are the eigenvalues, and  $\phi_k$  are eigenfunctions that form an orthonormal basis with a unit norm. Using the orthogonal expansion allows re-writing each functional observation  $X_i(t)$  in the Karhunen-Loève representation

$$X_i(t) = \mu(t) + \sum_{k=1}^{\infty} \xi_{ik} \phi_k(t) \quad \text{where} \quad \xi_{ik} = \int (X_i(t) - \mu(t)) \phi_k(t) dt. \quad (22)$$

The  $\xi_{ik}$  are the Functional Principal Components (FPC)[**Chiou'2014**]. Like PCA Principal Components, FPCs can also be clustered using clustering algorithms to find groups in the data.

The Karhunen-Loève representation enables a few different things. When using estimated forms of the mean function, FPCs, and eigenfunctions which can be found a number of ways such as numerical integration, using a limited number of FPCs that explain most of the variation in the

data can be used to generate noise reduced versions of the data while keeping the parts of the process that explain the most variability[ramsay2007applied]. If  $L$  is the number of FPC that explain for example 90% of the variation in the data, then noise reduced observations can be made by

$$\hat{X}_i(t) = \hat{\mu}(t) + \sum_{k=1}^L \hat{\xi}_{ik} \hat{\phi}_k(t). \quad (23)$$

This same method can be used to impute missing parts of the function as well[Chiou2014].

## Appendices

### C C++ Code used

#### C.1 fit.cpp (Model Main)

```
#define TMB_LIB_INIT R_init_fit
#include <TMB.hpp>
#include "pnorm4.hpp"
#include "Fcorr.hpp"
#include "MVMIX.hpp"

//Convert a matrix to it's proportions by row
template<class Type>
matrix<Type> propM(matrix<Type> M){
    int A = M.cols();
    int Y = M.rows();

    matrix<Type> retM(Y,A);

    for(int y = 0;y < Y;y++){
        vector<Type> cy = M.row(y);
        Type tot = cy.sum();
        vector<Type> prop = cy/tot;
        retM.row(y) = prop;
    }

    return retM;
}

/*Calculate population weighted average F for some set of ages
 bAge is the beginning age to calculate average F over, eAge is the end Age.*/
template<class Type>
vector<Type> aveF(matrix<Type> F, matrix<Type> N, int bAge,int eAge){
    int Y = F.rows();
    int dist = (eAge-bAge)+1;
    matrix<Type> aveF = F.array()*N.array();
    //Get the block of the matrix starting at year 0 going to year Y, column of bAge
    //extending the number of columns between eAge and bAge and take the sum rowwise.
    vector<Type> aveF_XY = aveF.block(0,bAge,Y,dist).rowwise().sum();
    vector<Type> tn = N.block(0,bAge,Y,dist).rowwise().sum();
    aveF_XY = aveF_XY/tn;
    return aveF_XY;
}

/*Calculate CRLs from catch proportions*/
template<class Type>
matrix<Type> makeCRLs(matrix<Type> cProps){
    int A = cProps.cols();
```

```

int Y = cProps.rows();

matrix<Type> crls(Y,A-1);

for(int a = 0; a < A-1;a++){
    for(int y = 0; y < Y;y++){
        vector<Type> cprow = cProps.row(y);
        //get p_a+...+p_A
        Type denom = cprow.segment(a,A-a).sum();
        Type num = cProps(y,a);
        Type pi = num/denom;
        crls(y,a) = log(pi/(1-pi));
    }
}

return crls;
}

/*Beverton-Holt Recruitment function*/
template<class Type>
Type bHolt(Type SSBly,vector<Type> bhParm){
    Type ret = bhParm(0) + log(SSBly) -log(1.0+exp(bhParm(1))*SSBly);
    return ret;
}

/*Ricker recruitment function*/
template<class Type>
Type ricker(Type SSBly,vector<Type> rickParm){
    Type ret = rickParm(0) + log(SSBly) - exp(rickParm(1))*SSBly;
    return ret;
}

// ' Smooth Hockey Stick Recruitment Function
template<class Type>
Type sHS(Type SSBly,vector<Type> segParm,Type gammaSq){
    //Alpha and delta must be positive.
    Type delta = exp(segParm(1));
    Type alpha = exp(segParm(0));

    Type gm = gammaSq/4.0;
    Type sq1 = sqrt(delta*delta + gm);
    Type sq2 = sqrt((SSBly-delta)*(SSBly-delta) +gm);
    Type ret = log(alpha*(SSBly+sq1-sq2));
    return ret;
}

template<class Type>
Type objective_function<Type>::operator() ()
{
    enum fleetType {
        survey = 0,
        Catch = 1,
        land = 2
    };

    enum recType {
        rw = 0,
        bh = 1,
        rick = 2,
        smoothHS = 3
    };

    // input data;
    DATA_MATRIX(M);
    DATA_MATRIX(weight);
    DATA_MATRIX(mat);

```

```

DATA_MATRIX(midy_weight);
DATA_VECTOR(log_index);
DATA_IVECTOR(i_zero); //A lot of this could be put in one IMATRIX?
DATA_IVECTOR(iyear);
DATA_IVECTOR(iage);
DATA_IVECTOR(isurvey);
DATA_IVECTOR(iq);
DATA_VECTOR(fs);
DATA_IVECTOR(ft);
DATA_VECTOR(lowerMult); //Lower multiplier for censored bounds
DATA_VECTOR(upperMult); //Upper multiplier for censored bounds
DATA_IVECTOR(keyF);
DATA_INTEGER(recflag);
DATA_INTEGER(corflag);
DATA_INTEGER(corflagCRL);
DATA_SCALAR(gammaSq); //For smoothHS recruitment function, controls smoothness of
    breakpoint
DATA_VECTOR(crlVec); //Vector of Continuation Ratio Logits for catch
DATA_IVECTOR(aveFages); //Ages in model ages to calculate average F over...0 bAge, 1 for
    eAge

DATA_VECTOR_INDICATOR(keep,log_index); //Used for one step predict, does not need to be
    read in from R
DATA_VECTOR_INDICATOR(keepCRL,crlVec);

int n = log_index.size();

int A = mat.cols();
int Y = mat.rows();
vector<Type> index = exp(log_index);

//define parameters;
PARAMETER(log_std_log_R);
PARAMETER_VECTOR(log_qparm);
PARAMETER_VECTOR(log_std_index);
PARAMETER_VECTOR(log_std_logF);
PARAMETER_VECTOR(log_std_CRL);
PARAMETER_VECTOR(log_std_landings); //std. dev for landings, size 0 if fit_land=0, else 1
PARAMETER(log_sdS);
PARAMETER_VECTOR(rec_parm);
PARAMETER_VECTOR(tRhoF);
PARAMETER_VECTOR(tRhoCRL);

//Random Effects
PARAMETER_MATRIX(log_F);
PARAMETER_MATRIX(log_N);

Type std_log_R = exp(log_std_log_R);
vector<Type> std_index = exp(log_std_index);
vector<Type> std_logF = exp(log_std_logF);
vector<Type> std_CRL = exp(log_std_CRL);
vector<Type> std_landings = exp(log_std_landings);
Type sdS = exp(log_sdS);

matrix<Type> F(Y,A);
matrix<Type> Z(Y,A);
matrix<Type> EC(Y,A);
matrix<Type> ECW(Y,A);

vector<Type> Elog_index(n);
vector<Type> std_index_vec(n);

//***** start the engine *****;

using namespace density;
Type nll = 0.0;

```

```

//compute F & Z
for(int a = 0;a < A;++a){
    for(int y = 0;y < Y;++y){
        //keyF really cuts down on number of random parameters
        F(y,a) = exp(log_F(y,keyF(a)));
        Z(y,a) = F(y,a) + M(y,a);
    }
}

matrix<Type> N = log_N.array().exp();

//Because of the magic from making log_N a parameter you can calculate SSB whenever you
want.
matrix<Type> B_matrix = weight.array()*N.array();
matrix<Type> SSB_matrix = mat.array()*B_matrix.array();
vector<Type> biomass = B_matrix.rowwise().sum();
vector<Type> ssb = SSB_matrix.rowwise().sum();
vector<Type> log_biomass = log(biomass);
vector<Type> log_ssb = log(ssb);

//Recruitment, adding recruitment curves much easier because ssb is already available and
will be optimized later
Type predN;
for(int y =1;y < Y;y++){
    switch(recflag){
        case rw:
            predN = log_N(y-1,0);
            break;
        case bh:
            predN = bHolt(ssb(y),rec_parm);
            break;
        case rick:
            predN = ricker(ssb(y),rec_parm);
            break;
        case smoothHS:
            predN = sHS(ssb(y),rec_parm,gammaSq);
            break;
        default:
            error("Not_□right_□rec_□type");
            break;
    }

    nll -= dnorm(log_N(y,0),predN,std_log_R,true);
}

//Fill in N, add survival process
for(int y = 1;y < Y;y++){
    for(int a = 1;a < A;++a){
        predN = log_N(y-1,a-1)-Z(y-1,a-1);
        if(a == A-1){//Plus group
            predN = logspace_add(predN,log_N(y-1,a)-Z(y-1,a));
        }
        nll -= dnorm(log_N(y,a),predN,sdS,true);
    }
}

//F correlation matrix making, sigmaFgen lives in Fcorr.hpp
matrix<Type> sigmaF = sigmaFgen(log_F.cols(),corflag,log_std_logF,tRhoF);
REPORT(sigmaF);

//Add the F random walk to the likelihood.
density::MVNORM_t<Type> Fdens(sigmaF);
for(int y = 1; y < Y;y++){
    nll += Fdens(log_F.row(y)-log_F.row(y-1));
}

```

```

//Expected catch, for making expected catch CRLs
//Expected landings too
for(int y = 0; y < Y; y++){
    for(int a = 0; a < A; a++){
        EC(y,a) = N(y,a)*(1.0-exp(-1.0*Z(y,a)))*F(y,a)/Z(y,a);
        ECW(y,a) = EC(y,a)*midy_weight(y,a);
    }
}

vector<Type> landings_pred = ECW.rowwise().sum();
vector<Type> total_expected_catch = EC.rowwise().sum();
matrix<Type> log_EC = EC.array().log();

vector<Type> log_landings_pred = log(landings_pred);
vector<Type> log_total_expected_catch = log(total_expected_catch);

REPORT(log_landings_pred);
REPORT(log_total_expected_catch);

// Landings, Survey index predictions
int ia, iy;
for(int i = 0; i < n; ++i){
    int fType = ft(i);
    ia = iage(i);
    iy = iyear(i);

    //Observation equations for survey and landings
    switch(fType){
    case survey:
        std_index_vec(i) = std_index(isurvey(i));
        Elog_index(i) = log_qparm(iq(i)) + log_N(iy,ia) - fs(i)*Z(iy,ia);
        if(i_zero(i) == 1){ //if zero use left-censor bound
nll -= keep(i)*pnorm4(log_index(i), Elog_index(i), std_index_vec(i), true);
        }
        else{
nll -= keep(i)*dnorm(log_index(i), Elog_index(i), std_index_vec(i), true);
        }
        break;
    case land:
        Elog_index(i) = log_landings_pred(iy);
        std_index_vec(i) = std_landings(0);
        nll -=
            keep(i)*censored_bounds(log_index(i), Elog_index(i), std_index_vec(i), -log(lowerMult(iy)), log(upperMult(iy)));
        break;
    default:
        error("Fleet_Type_not_implemented");
        break;
    }
}

matrix<Type> EPC = propM(EC); //Expected Catch Proportions
matrix<Type> ECRL = makeCRLs(EPC); //Expected Catch Continuation Ratio Logits

matrix<Type> S = propM(F); //Selectivity  $F_{\{y,a\}}/(\sum_{\{a\}}(F_{\{y,a\}}))$ 

REPORT(S);
REPORT(EPC);
REPORT(ECRL);

//Correlation matrix for CRLs, might as well reuse sigmaFgen...
matrix<Type> sigmaCRL = sigmaFgen(ECRL.cols(), corflagCRL, log_std_CRL, tRhoCRL);
REPORT(sigmaCRL);

//MV NORMAL MIXTURE distribution, second argument is the mixture probability, grabbed
    from SAM lives in MVMIX.hpp

```

```

//MVNORM_t in TMB doesn't support OSA residuals yet, MVMIX_t from SAM does however there
// is a slight speed loss
MVMIX_t<Type> CRLdens(sigmaCRL,0); //cause this is a MV MIXTURE duh. I could maybe
// actually try to use this part?
for(int y = 0;y < Y;y++){
  //IT'S VERY VERY IMPORTANT CATCH CRLs ARE SORTED YEAR THEN AGE HERE.
  vector<Type> ECRLr = ECRL.row(y);
  nll += CRLdens(crlVec.segment(y*(A-1),(A-1))-ECRLr,keepCRL.segment(y*(A-1),(A-1)));
}

//pop size weighted ave F;

vector<Type> aveFXY = aveF(F,N,aveFages(0),aveFages(1));
vector<Type> log_aveFXY = log(aveFXY);

REPORT(std_logF);
REPORT(std_log_R);
REPORT(std_index);

REPORT(N);
REPORT(F);
REPORT(Z);
REPORT(B_matrix);
REPORT(SSB_matrix);
REPORT(biomass);
REPORT(ssb);
REPORT(aveFXY);

REPORT(EC);
REPORT(ECW);

REPORT(landings_pred);

REPORT(Elog_index);

REPORT(log_F);

REPORT(log_std_index);
REPORT(log_qparm);

ADREPORT(log_landings_pred);
ADREPORT(log_biomass);
ADREPORT(log_ssb);
ADREPORT(log_aveFXY);
ADREPORT(log_qparm);

ADREPORT(std_logF);
ADREPORT(std_log_R);
ADREPORT(std_index);

return nll;
}

```

## C.2 pnorm4.hpp

```

//Log of dnorm
template<class Type>
Type ldnorm(Type x){
  Type ret = 0.5*(-x*x-log(2.0*M_PI));
  return ret;
}

//Atomic function for log(pnorm)
//Needed to get good stable derivatives for use with pnorm4
//Atomic double specifies the forward mode derivative, which is just pnorm(x,log.p=TRUE) =
  log(pnorm)

```

```

//Reverse is the explicit derivative of log(pnorm) = dnorm/pnorm =
    exp(log(dnorm)-log(pnorm)), more stable derivatives at
//extreme values.
TMB_ATOMIC_VECTOR_FUNCTION(
    // ATOMIC_NAME
    lpnorm1
    ,
    // OUTPUT_DIM
    1
    ,
    // ATOMIC_DOUBLE
    ty[0] = atomic::Rmath::Rf_pnorm5(tx[0],0,1,1,1);
    ,
    // ATOMIC_REVERSE
    Type pn = ty[0]; //Get the log-supporting pnorm
    px[0] = exp(ldnorm(tx[0])-pn)*py[0]; //derivative of log(pnorm) in a more
        computationally stable way
)

//Type function to link to atomic
template<class Type>
Type lpnorm1(Type x){
    CppAD::vector<Type> tx(1);
    tx[0] = x;
    return lpnorm1(tx)[0];
}

//vectorize it
VECTORIZED1_t(lpnorm1);

//A version of pnorm taking 4 arguments,
//q quantile of the normal distribution
//mean the mean of the normal distribution
//sd the standard deviation
//give_log if 1, returns log(pnorm(q,mean,sd))
template<class Type>
Type pnorm4(Type q, Type mean = 0., Type sd = 1., int give_log = 1){
    if(give_log == 0){
        return pnorm(q,mean,sd);
    }else{
        Type m = (q-mean)/sd;
        return lpnorm1(m);
    }
}

VECTORIZED1_t(pnorm4);
VECTORIZED4_ttti(pnorm4);

//Version for censored upper and lower bounds

//logspace_sub to do log(exp(log(x))-exp(log(y))) safely
extern "C" {
    /* See 'R-API: entry points to C-code' (Writing R-extensions) */
    double Rf_logspace_sub (double logx, double logy);
}

//Here x is std. normal
//This is the forward double function needed by the atomic function
//Returns log(exp(log(pnorm(x+upper)))-exp(log(pnorm(x-lower)))) but in a much more
    computationally stable way.
double censored_bou(double x,double lowerb,double upperb){
    double upper,lower;
    if(x > 0){
        lower = atomic::Rmath::Rf_pnorm5(-(x+upperb),0,1,1,1);
        upper = atomic::Rmath::Rf_pnorm5(-(x-lowerb),0,1,1,1);
    }else{
        upper = atomic::Rmath::Rf_pnorm5(x+upperb,0,1,1,1);

```



```

    //Rcout << "Upper bound: " << upper << std::endl;
    lower = atomic::Rmath::Rf_pnorm5(x-lowerb,0,1,1,1);
    //Rcout << "Lower bound: " << lower << std::endl;
}
double ret = Rf_logspace_sub(upper,lower);
//Rcout << "What we ret: " << ret << std::endl;
return ret;
}

TMB_ATOMIC_VECTOR_FUNCTION(
    // ATOMIC_NAME
    censored_b
    ,
    // OUTPUT_DIM
    1
    ,
    // ATOMIC_DOUBLE
    ty[0] = censored_bou(tx[0],tx[1],tx[2]);
    ,
    // ATOMIC_REVERSE
    Type cb = ty[0]; //Get the censored bounds from the forward mode.
    Type d1 = ldnorm(tx[0]+tx[2]);
    Type d2 = ldnorm(tx[0]-tx[1]);
    Type f1 = d1-cb;
    //Rcout << "This is f1: " << f1 << std::endl;
    Type f2 = d2-cb;
    //Rcout << "This is f2: " << f2 << std::endl;
    px[0] = (exp(f1)-exp(f2))*py[0];
)

template<class Type>
Type censored_b(Type x,Type lowerb,Type upperb){
    CppAD::vector<Type> tx(3);
    tx[0] = x;
    tx[1] = lowerb;
    tx[2] = upperb;
    return censored_b(tx)[0];
}

//Function to calculate log(pnorm(x+upper,mu,sd)-pnorm(x-lower,mu,sd)) with stable
derivatives
template<class Type>
Type censored_bounds(Type x,Type mu,Type sd,Type lower,Type upper){
    Type z = (x-mu)/sd;
    Type upperb = upper/sd;
    Type lowerb = lower/sd;
    return censored_b(z,lowerb,upperb);
}

```

### C.3 Fcorr.hpp

```

template<class Type>
Type rhoTrans(Type x){
    return Type(2)/(Type(1) + exp(-Type(2)*x))-Type(1);
}

VECTORIZED1_t(rhoTrans);

/*
    Generate correlation matrix for F's

    param nage number of ages
    param corrFlag Type of correlation structure
    param tRhoF vector of untransformed rhos for corr structure
*/

```

```

template<class Type>
matrix<Type> sigmaFgen(int nage,int corrFlag, vector<Type> logsdF,vector<Type> tRhoF){

    enum corrType{
        independent = 0,
        parallel = 1,
        compound = 2,
        ARone = 3,
        custom = 4
    };

    matrix<Type> sigma(nage,nage);
    vector<Type> sdF = exp(logsdF);
    vector<Type> rhoF = rhoTrans(tRhoF);

    sigma.setZero();

    switch(corrFlag){
    case independent:
        sigma.diagonal() = sdF*sdF;
        break;
    case parallel:
        sigma.diagonal() = sdF*sdF;
        for(int i=0; i < nage; i++){
            for(int j=0; j < i; j++){
                //Anders explained this is needed to get parallel to coverage. VERY PICKY ON NUMBER OF
                //0.999
                sigma(i,j) = 0.999*sdF(i)*sdF(j);
                sigma(j,i) = sigma(i,j);
            }
        }
        break;
    case compound:
        sigma.diagonal() = sdF*sdF;
        for(int i = 0;i < nage; i++){
            for(int j = 0; j < i; j++){
                sigma(i,j) = rhoF(0)*sdF(i)*sdF(j);
                sigma(j,i) = sigma(i,j);
            }
        }
        break;
    case ARone:
        sigma.diagonal() = sdF*sdF;
        for(int i = 0;i < nage; i++){
            for(int j = 0; j < i; j++){
                sigma(i,j) = pow(rhoF(0),Type(i-j))*sdF(i)*sdF(j);
                sigma(j,i) = sigma(i,j);
            }
        }
        break;
    case custom:
        sigma.diagonal() = sdF*sdF;
        for(int i = 0;i < nage; i++){
            for(int j = 0; j < i; j++){
                sigma(i,j) = rhoF(i)*rhoF(j)*sdF(i)*sdF(j);
                sigma(j,i) = sigma(i,j);
            }
        }
        break;

    default:
        error("F_correlation_mode_not_supported");
        break;
    }

    return sigma;
}

```

## C.4 MVMIX.hpp

This code was taken from the SAM github page as it supports OSA residuals and the current MVNORM.t provided in TMB does not. The mixture probability was set to 0 in the model. It is also licensed as GPL-2 and thus fine for inclusion here.

```
/*Code ripped from SAM's SAM/stockassessment/inst/include/define.hpp github file.
   I'm not going to attempt writing my own MVNORM with support for OSA residuals for a class
   project...
   Maybe for fun later...*/

template<class Type>
Type logspace_add_p (Type logx, Type logy, Type p) {
    return log((Type(1)-p)*exp(logy-logx)+p)+logx; // the order of x and y is taylored for
    this application
}

template<class Type>
Type logdrobust(Type x, Type p){
    Type ld1=dnorm(x,Type(0.0),Type(1.0),true);
    if(p<Type(1.0e-16)){
        return ld1;
    }else{
        Type ld2=dt(x,Type(3),true);
        Type logres=logspace_add_p(ld2,ld1,p);
        return logres;
    }
}
VECTORIZATION2_tt(logdrobust)

template <class Type>
class MVMIX_t{
    Type halfLogDetS;
    Type p1; /*fraction t3*/
    matrix<Type> Sigma;
    vector<Type> sd;
    matrix<Type> L_Sigma;
    matrix<Type> inv_L_Sigma;
public:
    MVMIX_t(){}
    MVMIX_t(matrix<Type> Sigma_, Type p1_){
        setSigma(Sigma_);
        p1=p1_;
    }
    matrix<Type> cov(){return Sigma;}
    void setSigma(matrix<Type> Sigma_){
        Sigma = Sigma_;
        sd = sqrt(vector<Type>(Sigma.diagonal()));
        Eigen::LLT<Eigen::Matrix<Type,Eigen::Dynamic,Eigen::Dynamic> > llt(Sigma);
        L_Sigma = llt.matrixL();
        vector<Type> D=L_Sigma.diagonal();
        halfLogDetS = sum(log(D));
        inv_L_Sigma = atomic::matinv(L_Sigma);
    }
    void setSigma(matrix<Type> Sigma_, Type p1_){
        setSigma(Sigma_);
        p1=p1_;
    }
    /** \brief Evaluate the negative log density */
    Type operator()(vector<Type> x){
        vector<Type> z = inv_L_Sigma*x;
        return -sum(logdrobust(z,p1))+halfLogDetS;
    }
    Type operator()(vector<Type> x, vector<Type> keep){
        matrix<Type> S = Sigma;
        vector<Type> not_keep = Type(1.0) - keep;
        for(int i = 0; i < S.rows(); i++){
```

```

    for(int j = 0; j < S.cols(); j++){
      S(i,j) = S(i,j) * keep(i) * keep(j);
    }
    //S(i,i) += not_keep(i) * pow((Type(1)-p1)*sqrt(Type(0.5)/M_PI)+p1*(Type(1)/M_PI),2);
    //(t(1))
    S(i,i) += not_keep(i) *
      pow((Type(1)-p1)*sqrt(Type(0.5)/M_PI)+p1*(Type(2)/(M_PI*sqrt(Type(3))))),2);
  }
  return MVMIX_t<Type>(S,p1)(x * keep);
}

vector<Type> simulate() {
  int siz = Sigma.rows();
  vector<Type> x(siz);
  for(int i=0; i<siz; ++i){
    Type u = runif(0.0,1.0);
    if(u<p1){
      x(i) = rt(3.0);
    }else{
      x(i) = rnorm(0.0,1.0);
    }
  }
  x = L_Sigma*x;
  return x;
}
};

template <class Type>
MVMIX_t<Type> MVMIX(matrix<Type> Sigma, Type p1){
  return MVMIX_t<Type>(Sigma,p1);
}

```

## C.5 C++ code used outside model

I decided it would be faster to simply convert the code I wrote to do CRLs inside the model for use in R when I had to switch from doing all the CRL related stuff inside the model to both inside and out for the sake of OSA residuals. This is just a version of makeCRLs that uses the Eigen library directly and does not take the log for the CRL as that's done later inside R. Using the Rcpp package these can be used as functions inside R. They are quick but have zero checks.

```

#include <Rcpp.h>
#include <RcppEigen.h>

// [[Rcpp::depends(RcppEigen)]]

//' catchProp
//'
//' make a row prop
//'
//' Why doesn't this work
//'
//' @param catchy a row to proportionize
//' @export
// [[Rcpp::export]]
Eigen::VectorXd catchProp(Eigen::Map<Eigen::VectorXd > catchy){
  double tot = catchy.sum();
  Eigen::VectorXd prop = catchy/tot;
  return prop;
}

//' makeCRLs
//'
//' Turn matrix of proportions into CRLs
//'
//' @param cProps matrix of proportions

```

```

//' @export
// [[Rcpp::export]]
Eigen::MatrixXd makeCRLs(Eigen::Map<Eigen::MatrixXd> cProps){
  int A = cProps.cols();
  int Y = cProps.rows();

  Eigen::MatrixXd crls(Y,A-1);

  for(int a = 0;a < A-1;a++){
    for(int y = 0;y < Y;y++){
      Eigen::VectorXd cprow = cProps.row(y);
      double denom = cprow.segment(a,A-a).sum();
      double num = cProps(y,a);
      double pi = num/denom;
      crls(y,a) = pi/(1-pi);
    }
  }
  return crls;
}

```

## C.6 R Code used

### C.6.1 spam.R

```

#'Fit the SPAM model
#'
#' Basic wrapper so I get a class to make things like AIC easier.
#'
#' @param dat list of data and parameters from spamSetup
#'
#' @useDynLib 'fit'
#'
#' @export
#'
spamFit <-
  function(data,parameters,indices,cVec,random=NULL,map=list(),control=list(),silent=TRUE,sdrep=TRUE,...){

    obj <- TMB::MakeADFun(data,parameters,map=map,random=random,silent=silent,...)

    opt <- nlminb(obj$par,obj$fn,obj$gr,control=control)

    fit <- list()

    fit$dat <- data
    fit$parameters <- parameters
    fit$map = map
    fit$obj <- obj
    fit$opt <- opt
    fit$rep <- obj$report()
    if(sdrep == TRUE){
      fit$sdr <- TMB::sdreport(fit$obj)
    }

    fit$indices = indices
    fit$cVec = cVec

    class(fit) <- "spam"
    print(opt$message)
    fit
  }

#'Returns the loglikelihood of a spam object
#'

```

```

#' @param object the spam object to get the loglikelihood of
#' @export
logLik.spam <- function(object,...)
{
  val = -object$opt$objective
  attr(val,"df") = length(object$opt$par)
  attr(val,"nobs") = length(c(object$dat$log_index,object$dat$crlVec))
  class(val) = "logLik"
  val
}

#'Returns the one step ahead residuals of a spam object and some other stuff for plotting
#'
#' @param object the spam object
#' @param method the one step ahead residuals to use
#' @param ... options to give oneStepPredict
#'
#' @export
residuals.spam <- function(object,method="oneStepGaussianOffMode",...){
  ospIL =
    TMB::oneStepPredict(object$obj, "log_index", "keep", discrete=FALSE, method=method,...)
  ospCRL =
    TMB::oneStepPredict(object$obj, "crlVec", "keepCRL", discrete=FALSE, method=method,...)

  ##Add crud for plotting
  ospIL =
    cbind(ospIL, Age=object$indices$Age, survey=object$indices$survey, Year=object$indices$Year, EV=object$
      i_zero=object$indices$i_zero)
  ospCRL =
    cbind(ospCRL, Age=object$cVec$Age, survey=object$cVec$survey, Year=object$cVec$Year, EV=as.vector(t(ob
      i_zero=object$cVec$i_zero)
  ret = rbind(ospIL, ospCRL)
  ret
}

#'Helper function for generating plusGroup information
pGroup <- function(x, plusGroup, sOrM="survey", ages, years){
  x <- x[x$Age %in% ages,]
  x <- x[x$Year %in% years,]
  if(!is.null(plusGroup)){
    x <- dplyr::mutate(x, Age = ifelse(Age > plusGroup, plusGroup, Age))
  }
  if(sOrM == "survey"){
    x <- dplyr::group_by(x, Year, Age, survey, fs)
    if(!is.null(plusGroup))
      x <- dplyr::summarize(x, index = sum(index, na.rm=TRUE))
    x <- dplyr::ungroup(x)
    x <- dplyr::arrange(x, survey, fs, Year, Age)
    x <- as.data.frame(x)
  }
  else if(sOrM == "catch"){
    x <- dplyr::group_by(x, Year, Age)
    if(!is.null(plusGroup))
      x <- dplyr::summarize(x, index = sum(index, na.rm=TRUE))
    x <- dplyr::ungroup(x)
    x <- dplyr::arrange(x, Year, Age)
    x <- tidyr::spread(x, key=Age, value=index)
    x <- dplyr::select(x, -Year)
    x <- as.matrix(x)
  }
  else{
    x <- dplyr::group_by(x, Year, Age)
    if(!is.null(plusGroup))
      x <- dplyr::summarize(x, index = mean(index, na.rm=TRUE))
    x <- dplyr::ungroup(x)
    x <- dplyr::arrange(x, Year, Age)
  }
}

```

```

        x <- tidyr::spread(x, key=Age, value=index)
        x <- dplyr::select(x, -Year)
        x <- as.matrix(x)
      }

      x
    }

#'Setup data for SPAM model
#'
#'Constructs data and parameters to be usable by SPAM.
#' This accepts the data setup by mdata3Ps.R in the SPAM data-raw folder.
#'
#'
#' @param surveys named list of surveys/catch to include, catch must be named catch in list
#' @param landings data.frame of landings to include with column for upper and lower
#'   censored bound multipliers
#' @param stock_wt data.frame of stock weights
#' @param midy_wt data.frame of midyear weights
#' @param mat data.frame of maturity
#' @param recType recruitment option
#' @param ages optional vector of ages to include, if using plus group, keep ages above
#'   plusGroup in vector
#' @param years optional vector of years to include
#' @param plusGroup optional age to specify as a plus group, indices are summed by that age
#'   and weights taken to be the mean of ages in plus group
#' @param idetect The survey detection limit
#'
#' @importFrom Rcpp sourceCpp
#'
#' @useDynLib SSM3PS
#'
#' @export
#'
datSetup <-
function(surveys, landings, stock_wt, midy_wt, mat, M=0.2, ages=NULL, years=NULL, plusGroup=NULL, idetect=0.0005)
if(is.null(ages)){
  ageL <- grep("Age", names(mat))
  ages <- as.numeric(gsub("Age", "", names(mat)[ageL]))
}

if(is.null(years)){
  years = min(mat$Year):max(mat$Year)
}

if(length(M) == 1){
  if(!is.null(plusGroup)){
    maxAge = plusGroup
  }else{
    maxAge = max(ages)
  }
  M = matrix(M, nrow=length(years), ncol=length(min(ages):maxAge))
}else{
  stopifnot(nrow(M) != length(years) & ncol(M) != length(ages), "M is not right size matrix")
}

vMake <- function(x){
  ageL <- grep("Age", names(x))
  v <- tidyr::gather(x, key="Age", value="index", ageL)
  v$Age <- as.numeric(gsub("Age", "", v$Age))
  v
}

##Handle the surveys into the indices
surVec <- lapply(seq_along(surveys), function(x, nam, i){

```

```

    ageL <- grep("Age", names(x[[i]]))
    surV <- tidyr::gather(x[[i]], key="Age", value="index", ageL)
    surV$Age <- as.numeric(gsub("Age", "", surV$Age))
    surV$survey = nam[[i]]
    surV
  }, x=surveys, nam=names(surveys))
  surVec <- do.call(rbind, surVec)
  surVec <- pGroup(surVec, plusGroup, sOrM="survey", ages, years)
  surVec$iage <- surVec$Age - min(ages)
  surVec$iyear <- surVec$Year - min(years)
  surVec$i_zero <- ifelse(surVec$index < idetect | is.na(surVec$index), 1L, 0L)
  surVec$qname <- paste0(surVec$survey, ":", stringr::str_pad(surVec$Age, 2, pad="0",))
  surVec$qname <- ifelse(surVec$survey == "catch" | surVec$survey ==
    "landings", NA, surVec$qname)
  surVec$iq <- as.numeric(factor(surVec$qname, levels=sort(unique(surVec$qname))))-1
  surVec$ft <- dplyr::case_when(surVec$survey == "catch" ~ 1L,
    surVec$survey == "landings" ~ 2L,
    TRUE ~ 0L) ##Otherwise it's a survey
  surVec$isurvey <- ifelse(surVec$survey == "catch" | surVec$survey ==
    "landings", NA, surVec$survey)
  surVec$isurvey <- as.numeric(factor(surVec$isurvey))-1

  ##Setup for the plus group...
  matVec <- vMake(mat)
  stockVec <- vMake(stock_wt)
  midyVec <- vMake(midy_wt)

  pMat <- pGroup(matVec, plusGroup, sOrM="notS", ages, years)
  pStock <- pGroup(stockVec, plusGroup, sOrM="notS", ages, years)
  pMidy <- pGroup(midyVec, plusGroup, sOrM="notS", ages, years)

  landings$landings <- landings$landings/1000
  landings <- landings[landings$Year %in% years,]
  UB = landings$UB
  LB = landings$LB
  landings <- data.frame(Year=landings$Year, Age=rep(NA, nrow(landings)),
    survey="landings", fs=rep(NA, nrow(landings)),
    index = landings$landings, iage=rep(NA, nrow(landings)),
    iyear=landings$Year-min(years), i_zero=0,
    qname=rep(NA, nrow(landings)), iq=rep(NA, nrow(landings)),
    ft=2, isurvey=rep(NA, nrow(landings)))
  surVec <- rbind(surVec, landings)

  surVec$index[is.na(surVec$index)] = 0
  surVec$index[surVec$index < idetect] = idetect
  surVec$index[surVec$survey == "catch"] = surVec$index[surVec$survey == "catch"]/1000

  cVec = surVec[surVec$survey == "catch",]
  surVec = surVec[surVec$survey != "catch",]
  catt = cVec[, c("index", "Age", "Year")]
  catt = tidyr::spread(catt, "Age", "index")
  cprop = t(apply(catt[, -1], 1, catchProp))
  crls = makeCRLs(cprop)
  crls = as.data.frame(cbind(years, crls))
  a2 = sort(unique(cVec$Age))
  a2 = a2[1:(length(a2)-1)]
  names(crls) = c("Year", paste0("Age", a2))
  crls = tidyr::gather(crls, "Age", "index", -Year)
  crls$Age = as.numeric(gsub("Age", "", crls$Age))
  cVec$index = NULL
  crls = dplyr::left_join(crls, cVec)
  crls = dplyr::arrange(crls, iyear, iage) ##Sorted for my shitty hack
  ##surVec = rbind(crls, surVec)

  tmb.data <- list(
    M = M,
    weight = pStock,
    mat = pMat,

```



```

    midy_weight = pMidy,
    log_index = log(surVec$index),
    i_zero = surVec$i_zero,
    iyear = surVec$iyear,
    iage = surVec$iage,
    isurvey = surVec$isurvey,
    iq = surVec$iq,
    fs = surVec$fs,
    ft = surVec$ft,
    lowerMult = LB,
    upperMult = UB,
    keyF = rep(0, length(unique(na.omit(surVec$iage)))),
    recflag = 0,
    corflag = 0,
    corflagCRL = 0,
    gammaSq = 0.1,
    crlVec = log(crls$index),
    aveFages = aveFages
  )

  dat <- list()
  dat$data <- tmb.data
  dat$param <- paramSetup(tmb.data)
  dat$indices <- surVec
  dat$crls <- crls
  dat
}

#'Basic Setup of parameters
#'
#' @export
paramSetup <- function(dat){
  A = ncol(dat$mat)
  Y = nrow(dat$mat)

  parm <- list(
    log_std_log_R = 0,
    log_qparm=rep(0, length(unique(na.omit(dat$iq)))),
    log_std_index=rep(log(0.3), length(unique(na.omit(dat$isurvey)))),
    log_std_logF = rep(log(0.1), length(unique(dat$keyF))),
    log_std_CRL = rep(0, A-1),
    log_std_landings = log(0.02),
    log_sdS = 0,
    rec_parm = numeric(0),
    log_F=matrix(log(0.3), nrow=Y, ncol=length(unique(dat$keyF))),
    log_N = matrix(0, nrow=Y, ncol=A)
  )

  parm
}

```

## C.6.2 NCplots.R

```

#'Plot line graph with confidence interval from sd.report object
#'
#' Given a name of something in an sd.report object and the number of corresponding
#' years, plot the time series and the corresponding confidence intervals and return a
#' ggplot object
#'
#' @param sName the name of the series to plot
#' @param sd.rep the sd.report object
#' @param years vector of corresponding years
#' @param alpha alpha value for confidence intervals
#' @param exp take the exponent of the value?
#' @param add return a list of gproto objects to add together plots?
#' @param color the color of the line

```

```

#' @param fcolor the color of the confidence interval
#' @param relYear year for everything else to be relative to
#'
#' @export
#'
tsCIplot <-
  function(sName, sd.rep, years, alpha=0.05, exp=FALSE, add=FALSE, color="black", fcolor="grey70", relYear=NULL) {
    crit = qnorm(alpha/2, lower.tail=FALSE)

    ind = names(sd.rep$value) == sName
    stdi = sd.rep$sd[ind]
    if(exp == TRUE){
      low = exp(sd.rep$value[ind] - crit*stdi)
      val = exp(sd.rep$value[ind])
      high = exp(sd.rep$value[ind] + crit*stdi)
    } else {
      low = sd.rep$value[ind] - crit*stdi
      val = sd.rep$value[ind]
      high = sd.rep$value[ind] + crit*stdi
    }
    dat = data.frame(val=val, year=years, high=high, low=low)
    if(!is.null(relYear)){
      ryear = dat[dat$year == relYear,]
      dat$val = dat$val/ryear$val
      dat$high = dat$high/ryear$high
      dat$low = dat$low/ryear$low
    }

    pRet = ggplot2::ggplot() +
      ggplot2::geom_line(ggplot2::aes(x=year, y=val), dat, color=color) +
      ggplot2::geom_ribbon(ggplot2::aes(x=year, ymin=low, ymax=high), dat, fill=fcolor, alpha=0.3)

    if(add==TRUE){
      pRet = list(line = ggplot2::geom_line(ggplot2::aes(x=year, y=val), dat, color=color),
        ribbon =
          ggplot2::geom_ribbon(ggplot2::aes(x=year, ymin=low, ymax=high), dat, alpha=0.3, fill=fcolor)
      )
    }

    pRet
  }

#'bubble plot of residuals
#'
#'Bubble plots of the residuals using either oneStepPredict residuals or not
#'
#' @param rep the report object of the fit
#' @param indices the indices data frame used in the data
#' @param oneStep optional oneStepPredict residuals
#'
#' @export
bubblePlot <- function(rep, indices, oneStep=NULL, pch=21){
  dat = data.frame(Age = indices$Age, Year = indices$Year,
    Cohort = indices$Year-indices$Age,
    survey = indices$survey,
    pred = rep$Elog_index,
    index = indices$index,
    log_index = log(indices$index),
    resid = rep$resid_index,
    resid_std = rep$std_resid_index)

  if(!is.null(oneStep)){
    dat$resid_std = oneStep$residual
  }
}

```

```

dat$colr = "deepskyblue"
dat$colr[dat$resid > 0] = 'firebrick2'
dat$pch = pch
dat$colr[(indices$i_zero == 1) & (dat$resid > 0)] = 'azure2'
dat$colr[(indices$i_zero == 1) & (dat$resid <= 0)] = 'dimgray'

jDarkGray <- 'grey20'
pRet = lattice::xyplot(factor(Age) ~ Year | survey, data=dat,
  ylab="Age", main="Residuals",
  cex = 1.5*sqrt(abs(dat$resid_std)/pi), fill.color=dat$colr,
  col=jDarkGray,
  par.strip.text=list(cex=0.5),
  panel=function(x,y,...,cex,fill.color,pch,subscripts){
    panel.abline(v=seq(min(dat$Year),max(dat$Year),by=5),
      col.line=grey(0.9))
    panel.xyplot(x,y,cex=cex[subscripts],
      pch=dat$pch,fill=fill.color[subscripts],...)
  })
pRet
}

#' Plot estimated q patterns for surveys
#'
#' Function to plot the estimated q catchability patterns for surveys.
#' Returns a ggplot2 plot. Values are standardized to the max of the q.
#'
#' @param sd.rep the sd.rep object
#' @param indices the indices object
#' @param alpha alpha value for the confidence intervals
#' @param qpname name of q-parameters in report if different
#'
#' @export
qPattern <- function(sd.rep, indices, alpha=0.05, qpname="log_qparm"){
  crit = qnorm(alpha/2, lower.tail=FALSE)

  ind = names(sd.rep$value) == qpname

  indy = indices[!is.na(indices$i_q),] ## Needed because catch is now together with surveys
  qest = aggregate(indy$i_q, list(age=indy$Age, survey=indy$survey), unique)
  qest$i_q = qest$x+1
  temp = as.numeric(table(qest$i_q))
  qest$nx = temp[qest$i_q]

  qest$Eest = exp(sd.rep$value[ind]);
  qest$sd = sd.rep$sd[ind]
  qest$E.L = exp(sd.rep$value[ind] - crit*qest$sd);
  qest$E.U = exp(sd.rep$value[ind] + crit*qest$sd);

  qest$colr='grey'
  qest$colr[qest$nx>1]='black'

  temp = tapply(qest$Eest, qest$survey, max)
  qest$max_Eest = temp[match(qest$survey, names(temp))]

  qest$sEest = qest$Eest/qest$max_Eest
  qest$sE.L = qest$E.L/qest$max_Eest
  qest$sE.U = qest$E.U/qest$max_Eest

  qest$survey1 = paste(qest$survey, '□max□□□', round(qest$max_Eest, digits=3), '□
    (x10000)', sep='')

  qest = qest[order(qest$survey1, qest$age),]

  ind1 = qest$survey %in% names(temp)

  pRet = ggplot2::ggplot(qest, ggplot2::aes(x=age, y=sEest)) + ggplot2::geom_line() +

```

```

    ggplot2::geom_point() +
    ggplot2::geom_ribbon(ggplot2::aes(ymin=sE.L,ymax=sE.U),alpha=0.3) +
      ggplot2::facet_wrap(~survey1) + ggplot2::xlab("Age_Class") +
      ggplot2::ylab("Survey_Q_Pattern")

  pRet
}

#'Corrplot of elements
#'
#' Corrplot of elements, uses ggcorrplot because I prefer storable plots...
#'
#' @param sd.rep
#'
#' @export
#'
corrParm <- function(sd.rep){
  dsd = sqrt(diag(sd.rep$cov.fixed))
  corr.matrix = diag(1/dsd)%*%sd.rep$cov.fixed%*%diag(1/dsd)
  ##Patch for all the identical names
  dsdU = make.unique(names(dsd))
  rownames(corr.matrix) = dsdU
  colnames(corr.matrix) = dsdU

  p.mat = ggcorrplot::cor_pmat(corr.matrix)

  pRet =
    ggcorrplot::ggcorrplot(corr.matrix,method="circle",p.mat=p.mat,insig="blank",colors=c("red","white"))
  pRet
}

#' Plot a lot of different things
#'
#' Compare the predicted to actual value of the survey or catch crl index by age or not.
#' Returns a ggplot plot. This works for both catch and surveys because everything is done
  in one
  vector now.
#'
#' @param resObj residuals spam object with extra junk for plotting
#' @param indName optional the name of the catch,survey,landing etc. to look at
#' @param facetAges facet wrap the plots by age? Takes sum of index and mean of residuals
  over ages if false
#' @param residuals plot the residuals instead of the line plots
#' @param oneStep optional oneStepPredict residuals to supply
#' @param loess add a loess smoother to the residuals?
#' @param bubble bubble plots of residuals?
#'
#' @export
oneStopPlot <-
  function(resObj,indName=NULL,facetAges=TRUE,residuals=FALSE,loess=FALSE,bubble=FALSE){
    indy = resObj
    indy$Age = as.factor(indy$Age)
    indy$zero = factor(indy$i_zero,labels=c("nonzero","zero"),levels=c(0,1))

    if(!is.null(indName)){
      indy = indy[indy$survey %in% indName,]
    }

    indy$zero2 = dplyr::case_when((indy$i_zero==1)&(indy$residual >0) ~ "goodZero",
      (indy$i_zero==1)&(indy$residual <= 0) ~ "badZero",
      indy$residual > 0 ~ "pos",
      indy$residual <= 0 ~ "neg")

    if(facetAges == FALSE){
      indy = dplyr::group_by(indy,Year,survey)
      indy =

```

```

        summarize(indy, observation=sum(observation), EV=sum(EV), residual=mean(residual))
        facForm = as.formula(~ survey)
    }else{
        facForm = as.formula(~ survey + Age)
    }

    if(!("zero" %in% colnames(indy))){
        indy$zero = 0
        indy$zero = factor(indy$zero, labels=c("nonzero", "zero"), levels=c(0,1))
    }

    if(residuals == FALSE & bubble != TRUE){
        pRet = ggplot2::ggplot(indy, ggplot2::aes(x=Year, y=observation)) +
            ggplot2::geom_line(ggplot2::aes(x=Year, y=observation)) +
            ggplot2::geom_line(ggplot2::aes(x=Year, y=EV), color="red") +
            ggplot2::xlab("Year") + ggplot2::ylab("Log10 Index") +
            ggplot2::facet_wrap(facForm) +
            ggplot2::geom_point(ggplot2::aes(x=Year, y=observation, col=zero)) +
            ggplot2::scale_color_manual(values=c(zero="orange", nonzero="black"))
    }else if(residuals == TRUE && bubble == FALSE){
        pRet = ggplot2::ggplot(indy, ggplot2::aes(x=Year, y=residual)) +
            ggplot2::geom_point(ggplot2::aes(x=Year, y=residual, col=zero)) +
            ggplot2::xlab("Year") + ggplot2::ylab("Standardized10 Residuals") +
            ggplot2::facet_wrap(facForm) +
            ggplot2::scale_color_manual(values=c(zero="orange", nonzero="black"))
        if(loess==TRUE){
            pRet = pRet + ggplot2::geom_smooth()
        }
    }else if(bubble == TRUE){
        if(facetAges == FALSE){
            stop("This is a result of a bad design")
        }
        facForm = as.formula(~ survey)
        pRet = ggplot2::ggplot(indy, ggplot2::aes(x=Year, y=Age)) +
            ggplot2::geom_point(size=3*sqrt(abs(indy$residual)/pi))
        pRet = pRet +
            ggplot2::geom_point(ggplot2::aes(col=zero2, size=1.5*sqrt(abs(residual)/pi), stroke=0))
        pRet = pRet +
            ggplot2::scale_color_manual(values=c(pos="firebrick2", neg="deepskyblue", goodZero="green", badZero="red"))
        pRet = pRet + ggplot2::facet_wrap(facForm)
    }

    pRet
}

#' Residuals by year, cohort, age and against expected values
#'
#' Returns a list containing all four ggplot objects which can then easily be combined with
#' something like ggarrange.
#'
#' @param resObj the report object containing the residuals
#'
#' @export
#'
resid4plot <- function(dat){

    dat$cohort = dat$Year - dat$Age

    yplot = ggplot2::ggplot(dat, ggplot2::aes(x=Year, y=residual, label=as.character(Age))) +
        ggplot2::geom_text() + ggplot2::geom_smooth(color="red")

    cplot = ggplot2::ggplot(dat, ggplot2::aes(x=cohort, y=residual)) +
        ggplot2::geom_point(ggplot2::aes(cohort, residual, shape=3)) +

```

```

ggplot2::geom_smooth(color="red") + ggplot2::scale_shape_identity()

aplot = ggplot2::ggplot(dat,ggplot2::aes(x=Age,y=residual)) +
  ggplot2::geom_point(ggplot2::aes(Age,residual,shape=3)) +
  ggplot2::geom_smooth(color="red") + ggplot2::scale_shape_identity()

eplot = ggplot2::ggplot(dat,ggplot2::aes(x=EV,y=residual)) +
  ggplot2::geom_point(ggplot2::aes(x=EV,y=residual,shape=3)) +
  ggplot2::geom_smooth(color="red") + ggplot2::scale_shape_identity()

ret = list(yplot=yplot,cplot=cplot,aplot=aplot,eplot=eplot)
ret
}

```

### C.6.3 fpcaData.R

```

setwd("~/Documents/GitHub/SSM3PS/data-raw/data_EDA")
cnames = c('Year',paste0('Age',2:14))
DFO_stock_wt = read.table(file="DFO_stock_wt.txt",header=TRUE,col.names=cnames)
comm_wt = read.table(file='comm_wt.txt',header=TRUE)
cnames = c('Year',paste0('Age',1:16))
mat = read.table(file="mat.txt",header=FALSE,col.names=cnames)

##Weird old FDA tangent because
library(fdapace)

##LITERALLY DID NOT THINK ID BE USING THIS AGAIN
FPCAImpute <- function(FPCAobj,data,tau){
  L <- min(which(FPCAobj$cumFVE >= tau*100))
  impute <- matrix(0,nrow=nrow((FPCAobj$xiEst)),ncol=nrow(FPCAobj$phi))
  for(i in 1:length(FPCAobj$xiEst[,1])){
    impute[i,] = FPCAobj$xiEst[i,1:L]*t(FPCAobj$phi[,1:L]) + FPCAobj$mu
  }
  for(i in 1:length(data))
  {
    if(is.na(data[i])){
      data[i] <- impute[i]
    }
  }
  data
}

tcomm = t(comm_wt[comm_wt$Year %in% 1978:2016,-1])

ids = rep(1:ncol(tcomm),each=nrow(tcomm))
##Mid-year weights so everybody aged .5
s <- seq(3.5,14.5,by=1)
tv <- rep(s,ncol(tcomm))

fdaIn <- MakeFPCAInputs(IDs=ids,tVec=tv,tcomm)
cWClust <- FClust(fdaIn$Ly,fdaIn$Lt,k=3,optnsFPCA=list(FVEthreshold=0.9))
names(cWClust$cluster) = 1978:2014

tcomm = t(comm_wt[, -1])

ids = rep(1:ncol(tcomm),each=nrow(tcomm))
s <- seq(3.5,14.5,by=1)
tv <- rep(s,ncol(tcomm))

fdaIn <- MakeFPCAInputs(IDs=ids,tVec=tv,tcomm)
FPCAWeights <- FPCA(fdaIn$Ly,fdaIn$Lt,optns=list(dataType="Sparse"))

##Instead of ages 3.5 to 14.5, I need 1.5 to 16.5 for commercial weights and 1.0 to 16.0
  for stock
newGrid <- seq(0.5,16.5,by=1)
nMu <- spline(FPCAWeights$workGrid,FPCAWeights$mu,xout=newGrid)$y

```

```

nPhi <- apply(FPCAWeights$phi,2,function(y) spline(FPCAWeights$workGrid,y,xout=newGrid)$y)
nFPCA <- list(cumFVE = FPCAWeights$cumFVE,xiEst=FPCAWeights$xiEst,mu=nMu,phi=nPhi)

##Remove noise by extrapolating on the fitted values
ncomm = ConvertSupport(FPCAWeights$workGrid,s,phi=t(fitted(FPCAWeights,K=1)))
##Spots for missing ages
aM = rep(NA,ncol(ncomm))
ncomm = rbind(aM,ncomm)
ncomm = rbind(aM,ncomm)
ncomm = rbind(aM,ncomm)
ncomm = rbind(ncomm,aM)
ncomm = rbind(ncomm,aM)

fpcaComm_wt <- FPCAImpute(nFPCA,t(ncomm),0.9)

fComm_wt <- data.frame(Year=comm_wt$Year)
fComm_wt <- cbind(fComm_wt,fpcaComm_wt)
names(fComm_wt) = c("Year",paste0("Age",newGrid))
scm = tidyr::gather(fComm_wt,key="Age",value="Weight",-Year)
scm$Age = as.numeric(gsub("Age","",scm$Age))

fc2016 = apply(fComm_wt[fComm_wt$Year %in% 2012:2014,-1],2,median)
ffComm_wt = rbind(fComm_wt,c(2015,fc2016))
ffComm_wt = rbind(ffComm_wt,c(2016,fc2016))
ffComm_wt = ffComm_wt[,-2]
names(ffComm_wt) = c("Year",paste0("Age",1:16))
write.csv(ffComm_wt,"fComm_wt.csv")

##FakeSW are basically just a shift down the curve half a year
fakeSW = t(ConvertSupport(newGrid,1:16,phi=t(fpcaComm_wt)))
fakeSW = cbind(fComm_wt$Year,fakeSW)
fakeSW = as.data.frame(fakeSW)
names(fakeSW) = c("Year",paste0("Age",1:16))
fs2016 = apply(fakeSW[fakeSW$Year %in% 2012:2014,-1],2,median)
ffakeSW = rbind(fakeSW,c(2015,fs2016))
ffakeSW = rbind(ffakeSW,c(2016,fs2016))
write.csv(ffakeSW,"fStock_wt.csv")

ssw = tidyr::gather(fakeSW,"Age","Weight",-Year)
ssw$Age = as.numeric(gsub("Age","",ssw$Age))

#ggplot <- ggplot(scm,aes(Age,Weight)) + geom_point(aes(Age,Weight,group=Year),col="red") +
  geom_point(aes(Age,Weight,group=Year),data=ssw,col="blue")

##Stock weights are supposed to be January weights, extrapolate back to Jan.
##Using FPCA on DFO stock weights odd, so try back-calculating this way

##Try to impute ages 1,2,15,16 imputing at bounds is sort of gross, but the mean function
  is so simple this might actually
##be OKAY. Worth a shot anyways.

##I don't think 2 year old fish weigh nothing
tstock = t(DFO_stock_wt[DFO_stock_wt$Year %in% 1978:2019,-(1:2)])

ids = rep(1:ncol(tstock),each=nrow(tstock))
s <- 3:14
tv <- rep(s,ncol(tstock))

fdaIn <- MakeFPCAInputs(IDs=ids,tVec=tv,tstock)
DSClust <- FClust(fdaIn$Ly,fdaIn$Lt,optnsFPCA=list(FVEthreshold=0.9))
names(DSClust$cluster) = 1978:2017

tstock = t(DFO_stock_wt[,-(1:2)])

ids = rep(1:ncol(tstock),each=nrow(tstock))

```

```

s <- 3:14
tv <- rep(s, ncol(tstock))

fdaIn <- MakeFPCAInputs(IDs=ids, tVec=tv, tstock)
FPCASWeights <- FPCA(fdaIn$Ly, fdaIn$Lt, optns=list(dataType="Sparse"))

##Repeat steps for imputing ages
nstock = ConvertSupport(FPCASWeights$workGrid, 3:14, phi=t(fitted(FPCASWeights, K=1)))
##Spots for missing ages
aM = rep(NA, ncol(nstock))
nstock = rbind(aM, nstock)
nstock = rbind(aM, nstock)
nstock = rbind(nstock, aM)
nstock = rbind(nstock, aM)

nMu <- spline(FPCASWeights$workGrid, FPCASWeights$mu, xout=newGrid)$y
nPhi <- apply(FPCASWeights$phi, 2, function(y) spline(FPCASWeights$workGrid, y, xout=newGrid)$y)
nSFPCA <- list(cumFVE = FPCASWeights$cumFVE, xiEst=FPCASWeights$xiEst, mu=nMu, phi=nPhi)

fPCAStock_wt <- FPCAImpute(nSFPCA, t(nstock), 0.9)

fStock_wt <- data.frame(Year=DF0_stock_wt$Year)
fStock_wt <- cbind(fStock_wt, fPCAStock_wt)
names(fStock_wt) = c("Year", paste0("Age", 1:16))

##Now for Maturities!
##Squeeze to apply qlogis, so I can plogis on the return and ensure 0-1 bounds
mats = apply(mat[, -1], 1, function(x) mice::squeeze(x, bounds=c(0.0000000001, 0.9999999999)))
mats = qlogis(mats)

ids = rep(1:nrow(mats), each=ncol(mats))
s <- mat$Year
tv <- rep(s, nrow(mats))

fdaIn <- MakeFPCAInputs(IDs=ids, tVec=tv, t(mats))
FPCAMat <- FPCA(fdaIn$Ly, fdaIn$Lt, optns=list(methodMuCovEst="smooth", userBwMu=14))
fMaty <- fitted(FPCAMat, K=1)

fpMat <- plogis(fMaty)
fMat <- mat$Year
fMat <- cbind(fMat, t(fpMat))
fMat <- as.data.frame(fMat)

names(fMat) <- c("Year", paste0("Age", 1:16))
write.csv(fMat, file="fMat.csv")

```

## C.6.4 mdata3Ps.R

```

library(tidyverse)
library(lubridate)

setwd("../data-raw/data_EDA")

##landings
cnames = c('Year', 'landings')
landings = read.table(file='landings.dat', header=FALSE, col.names=cnames) %>%
  group_by(Year) %>%
  summarize(landings=sum(landings))

##Asked Danny Ings about bounds on landings broke it down into periods
##"Medium bounds"
##1959-1993, somewhat wide

```



```

landings$UB[landings$Year %in% 1959:1993] = 1.5
landings$LB[landings$Year %in% 1959:1993] = 0.5
## ##1994-96, bycatch only, narrow
landings$UB[landings$Year %in% 1994:1996] = 1.10
landings$LB[landings$Year %in% 1994:1996] = 0.90
## ##1997-98, maybe 99, very wide, wider than 1959-93
landings$UB[landings$Year %in% 1997:1999] = 1.75
landings$LB[landings$Year %in% 1997:1999] = 0.35
## ##1999(or 2000) to present, narrower than 1959-1993 but wider than 94-96
landings$UB[landings$Year %in% 2000:2016] = 1.25
landings$LB[landings$Year %in% 2000:2016] = 0.75

##Inshore/offshore survey
cnames = c('Year',paste0('Age',1:16))
RV_IO.matrix = read.table(file="RV_IO.dat",header=FALSE,col.names=cnames)
RV_IO.matrix$Year = as.character(RV_IO.matrix$Year)
yfs = do.call(rbind, strsplit(RV_IO.matrix$Year, "\\."))
RV_IO.matrix$Year = yfs[,1]
RV_IO.matrix$fs = as.numeric(yfs[,2])/10
RV_IO.matrix$Year = as.numeric(RV_IO.matrix$Year)
##Remove 2006
RV_IO.matrix = RV_IO.matrix[RV_IO.matrix$Year != "2006",]

##Offshore survey
cnames = c('Year',paste0('Age',1:16))
RV_OFF.matrix = read.table(file="RV_OFF.dat",header=FALSE,col.names=cnames)
RV_OFF.matrix$Year = as.character(RV_OFF.matrix$Year)
yfs = do.call(rbind, strsplit(RV_OFF.matrix$Year, "\\."))
RV_OFF.matrix$Year = yfs[,1]
RV_OFF.matrix$fs = as.numeric(yfs[,2])/10
RV_OFF.matrix$Year = as.numeric(RV_OFF.matrix$Year)
##Remove RV_OFF after 1996 since RV_IO is extension
RV_OFF.matrixR = RV_OFF.matrix[RV_OFF.matrix$Year < "1996",]

##Read in the GEAC survey data
cnames = c('Year', 'MedDate', paste0('Age', 1:16))
geac = read.csv("geac.csv", header=TRUE, col.names=cnames)
geac$fs <- lubridate::month(as.Date(geac$MedDate))/12 ##Yeah I know...
geac = geac[,-2]
##Remove 1997 and 2007
geacR = geac[geac$Year != "1997",]
geacR = geacR[geacR$Year != "2007",]

##Read in the French survey data
frenchS = read.csv("ERHAPS_indices_CAN.csv", sep=";")
frenchSur = data.frame(Year = 1978:1992)
frenchSur = cbind(frenchSur, t(frenchS)[-1,])
names(frenchSur) = c('Year', paste0('Age', 1:14))
rownames(frenchSur) = NULL

frenchDates = read.csv("ERHAPS_res.csv")
frenchDates$Start.date = as.Date(frenchDates$Start.date)
frenchDates$End.date = as.Date(frenchDates$End.date)
frenchDates$midDate = frenchDates$Start.date + ((frenchDates$End.date -
  frenchDates$Start.date)/2)
frenchDates$fs = lubridate::month(frenchDates$midDate)/12
frenchDates = frenchDates[,c("Year", "fs")]
frenchSur = dplyr::inner_join(frenchSur, frenchDates)

cnames = c('Year', paste0('Age', 1:16))
mat = read.table(file="mat.txt", header=FALSE, col.names=cnames)

midy_wt = read.table(file='midy_wt.dat', header=TRUE)

comm_wt = read.table(file='comm_wt.txt', header=TRUE)

cnames = c('Year', paste0('Age', 2:14))

```

```

DFO_stock_wt = read.table(file="DFO_stock_wt.txt",header=TRUE,col.names=cnames)

stock_wt = read.table(file='stock_wt.dat',header=TRUE)

cnames = c('Year',paste0('Age',2:16))
catch = read.table(file='catch.dat',header=FALSE,col.names=cnames)

catch$fs = NA
surveys = list(RVOff = RV_OFF.matrixR,RVInOff =
  RV_IO.matrix,geac=geacR,catch=catch,french=frenchSur)

##Read in FPCA files
fstock_wt = read.csv("fStock_wt.csv",header=TRUE)
fstock_wt = fstock_wt[,-1]
fcomm_wt = read.csv("fComm_wt.csv",header=TRUE)
fcomm_wt[,-1] = fcomm_wt[,-1]
fmat = read.csv("fMat.csv",header=TRUE)
fmat[,-1] = fmat[,-1]

rDat = list(surveys=surveys,fstock_wt=fstock_wt,fcomm_wt=fcomm_wt,fmat=fmat,
  catch=catch,DFO_stock_wt=DFO_stock_wt,stock_wt=stock_wt,midy_wt=midy_wt,mat=mat,
  RVOff=RV_OFF.matrix,RVInOff=RV_IO.matrix,geac=geac,french=frenchSur,landings=landings)
usethis::use_data(rDat,overwrite=TRUE)

```