

# Coverpage

```
## Document build date: Tue Dec 20 14:52:14 2016
## Working directory :
##       D:/projects/flr/drafting-doc
## Current contents of .GlobalEnv:
##       .First .Last
##
## Session information:
## R version 3.3.2 (2016-10-31)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows 10 x64 (build 10586)
##
## locale:
## [1] LC_COLLATE=English_United Kingdom.1252
## [2] LC_CTYPE=English_United Kingdom.1252
## [3] LC_MONETARY=English_United Kingdom.1252
## [4] LC_NUMERIC=C
## [5] LC_TIME=English_United Kingdom.1252
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] xtable_1.8-2      lattice_0.20-34    RColorBrewer_1.1-2
## [4] knitr_1.15.1
##
## loaded via a namespace (and not attached):
## [1] magrittr_1.5  tools_3.3.2  stringi_1.1.2 grid_3.3.2   stringr_1.1.0
## [6] evaluate_0.10
```

# Description of Catch at Age model

Colin P. Millar, Ernesto Jardim

## Abstract

This document presents the statistical catch-at-age stock assessment model developed in the JRC Assessment For All (**a4a**) initiative. The stock assessment model framework is a non-linear catch-at-age model implemented in R (<http://www.r-project.org/>) / FLR (<http://www.flr-project.org/>) / ADMB (<http://www.admb-project.org/>) that can be applied rapidly to a wide range of situations with low parametrization requirements. The model structure is defined by submodels, which are the different parts that require structural assumptions. There are 5 submodels in operation: a model for F-at-age, a model for the initial age structure, a model for recruitment, a (list) of model(s) for abundance indices catchability-at-age, and a list of models for the observation variance of catch-at-age and abundance indices. The submodels form use linear models. This opens the possibility of using the linear modelling tools available in R: see for example the mgcv (<http://cran.r-project.org/web/packages/mgcv/index.html>) gam formulas, or factorial design formulas using. Detailed model formulas, several diagnostic tools and a large set of models are presented in the document. Additionally, advanced features like external weighting of the likelihood components and MCMC fits are also described. The target audience for this document are readers with some experience in R and some background on stock assessment. The document explains the approach being developed by a4a for fish stock assessment and scientific advice. It presents a mixture of text and code, where the first explains the concepts behind the methods, while the last shows how these can be run with the software provided.

# Contents

<b>1</b>	<b>Background</b>	<b>4</b>
<b>2</b>	<b>Motivation</b>	<b>4</b>
<b>3</b>	<b>Data and model</b>	<b>4</b>
<b>4</b>	<b>Stock assessment model details</b>	<b>5</b>
4.1	Model fitting . . . . .	6
<b>5</b>	<b>Extending the model (brief)</b>	<b>7</b>
<b>6</b>	<b>Summary</b>	<b>7</b>
<b>7</b>	<b>Implementation (sketch)</b>	<b>8</b>

# 1 Background

The stock assessment model framework is a non-linear catch-at-age model implemented in R/FLR/ADMB that can be applied rapidly to a wide range of situations with low parametrization requirements.

In the **a4a** assessment model, the model structure is defined by submodels, which are the different parts of a statistical catch at age model that require structural assumptions.

There are 5 submodels in operation:

- a model for F-at-age,
- a (list) of model(s) for abundance indices catchability-at-age,
- a model for recruitment,
- a list of models for the observation variance of catch-at-age and abundance indices,
- a model for the initial age structure,

In practice, we fix the variance models and the initial age structure models, but in theory these can be changed.

The submodels form use linear models. This opens the possibility of using the linear modelling tools available in R: see for example the [mgcv](#) gam formulas, or factorial design formulas using `lm()`. In R's linear modelling language, a constant model is coded as  $\sim 1$ , while a slope over age would simply be  $\sim age$ . For example, we can write a traditional year/age separable F model like  $\sim factor(age) + factor(year)$ .

The 'language' of linear models has been developing within the statistical community for many years, and constitutes an elegant way of defining models without going through the complexity of mathematical representations. This approach makes it also easier to communicate among scientists

- 1965 J. A. Nelder, notation for randomized block design
- 1973 Wilkinson and Rodgers, symbolic description for factorial designs
- 1990 Hastie and Tibshirani, introduced notation for smoothers
- 1991 Chambers and Hastie, further developed for use in S

There are two basic types of assessments available in **a4a** : the management procedure fit and the full assessment fit. The management procedure fit does not compute estimates of covariances and is therefore quicker to execute, while the full assessment fit returns parameter estimates and their covariances at the expense of longer fitting time.

## 2 Motivation

The goal is to describe an age based model that is robust and easy to use. Robustness here means robust parameter estimation in addition to being robust to (underlying) model complexity. We aim to provide robust parameter estimates by concentrating parameters out of the objective function where possible and using structured random effects to allow for complexity when the data support it. More on this in the coming sections, but first a brief description of the data that the model will have to fit to.

## 3 Data and model

The data are

$C_{at}$  catch at age  $a$  and year  $t$

$S_{atk}$  abundance index for age  $a$  and year  $t$  from the  $k$ th survey or CPUE series,  $k = 1, 2, \dots$

The model is an age structure model where the number of fish in a given cohort  $N$  at the start of the following year is the number of fish that survived the perils of the current year. We assume that fish die through the year at a constant rate  $e^{-Z}$  ( $Z$  is positive), and that this rate is solely due to natural causes ( $M$ ) and fishing ( $F$ ) so that the total mortality rate is  $Z = F + M$ . This results in the model

$$N_{a+1,t+1} = N_{at}e^{-Z_{at}}$$

Abundance indices are observations of the relative abundance not of absolute abundance. This is because trawl surveys do not detect every fish but a fixed proportion  $Q$ . This proportion depends on age through length and means the index is proportional to abundance

$$S_{at} = Q_a N_{at}$$

If  $F$  and  $M$  are constant through the year catches arise as a fraction of those fish that died, and is written here as the familiar Baranov catch equation

$$\begin{aligned} C_{at} &= \frac{F_{at}}{Z_{at}} (N_{at} - N_{a+1,t+1}) \\ &= \frac{F_{at}}{Z_{at}} \left(1 - e^{-Z_{at}}\right) N_{at} \end{aligned}$$

These last two equations show that in there own way, catches and abundance indices are both observations of the numbers of fish in the population. Neither is sufficient to estimate the absolute abundances  $N$  but together they can be used to estimate both  $N$  and  $F$ . One way of doing this is using a statistical catch at age approach

## 4 Stock assessment model details

Modelled catches  $C$  are defined in terms of the three quantities, natural mortality  $M$ , fishing mortality  $F$  and recruitment  $R$ , using a modified form of the well known Baranov catch equation:

$$C_{ay} = \frac{F_{ay}}{F_{ay} + M_{ay}} \left(1 - e^{-(F_{ay} + M_{ay})}\right) R_y e^{-\sum (F_{ay} + M_{ay})}$$

where  $a$  and  $y$  denote age and year. Modelled survey indices  $I$  are defined in terms of the same three quantities with the addition of survey catchability  $Q$ :

$$I_{ays} = Q_{ays} R_y e^{-\sum (F_{ay} + M_{ay})}$$

where  $s$  denotes survey or abundance index and allows for multiple surveys to be considered. Observed catches  $C^{(obs)}$  and the observed survey indices  $I^{(obs)}$  are assumed to be log-normally distributed, or equivalently, normally distributed on the log-scale, with age, year and survey specific observation variance:

$$\log C_{ay}^{(obs)} \sim \text{Normal}\left(\log C_{ay}, \sigma_{ay}^2\right) \quad \log I_{ays}^{(obs)} \sim \text{Normal}\left(\log I_{ays}, \tau_{ays}^2\right)$$

The full log-likelihood for the **a4a** statistical catch at age model can now be defined as the sum of the log-likelihood of the observed catches ( $\ell_N$  is the log-likelihood of a normal distribution)

$$\ell_C = \sum_{ay} w_{ay}^{(c)} \ell_N\left(\log C_{ay}, \sigma_{ay}^2; \log C_{ay}^{(obs)}\right)$$

and the log-likelihood of the observed survey indices

$$\ell_I = \sum_s \sum_{ay} w_{ays}^{(s)} \ell_N\left(\log I_{ays}, \tau_{ays}^2; \log I_{ays}^{(obs)}\right)$$

giving the total log-likelihood

$$\ell = \ell_C + \ell_I$$

which is defined in terms of the strictly positive quantites,  $M_{ay}$ ,  $F_{ay}$ ,  $Q_{ays}$  and  $R_y$ , and the observation variances  $\sigma_{ay}$  and  $\tau_{ays}$ . As such, the log-likelihood is over-parameterised as there are many more parameters than observations. In order to reduce the number of parameters,  $M_{ay}$  is assumed known (as is common), and the remaining parameters are written in terms of a linear combination of covariates  $x_{ayk}$ , e.g.

$$\log F_{ay} = \sum_k \beta_k x_{ayk}$$

where  $k$  is the number of parameters to be estimated and is sufficiently small. Using this technique the quantities  $\log F$ ,  $\log Q$ ,  $\log \sigma$  and  $\log \tau$  (in bold in the equations above) can be described by a reduced number of parameters. The following section has more discussion on the use of linear models in **a4a**.

### Stock recruitment relationships

The **a4a** statistical catch at age model can additionally allow for a functional relationship to be imposed that links predicted recruitment  $\tilde{R}$  based on spawning stock biomass and modelled recruitment  $R$ , included as a fixed variance random effect. Options for the relationship are the hard coded models Ricker, Beverton Holt, smooth hockeystick or geometric mean. This is implemented by including a third component in the log-likelihood

$$\ell_{SR} = \sum_y \ell_N \left( \log \tilde{R}_y(a, b), \phi_y^2; \log R_y \right)$$

giving the total log-likelihood

$$\ell = \ell_C + \ell_I + \ell_{SR}$$

Using the (time varying) Ricker model as an example, predicted recruitment is

$$\tilde{R}_y(a_y, b_y) = a_y S_{y-1} e^{-b_y S_{y-1}}$$

where  $S$  is spawning stock biomass derived from the model parameters  $F$  and  $R$ , and the fixed quantites  $M$  and mean weights by year and age. It is assumed that  $R$  is log-normally distributed, or equivalently, normally distributed on the log-scale about the (log) recruitment predicted by the SR model  $\tilde{R}$ , with known variance  $\phi^2$ , i.e.

$$\log R_y \sim \text{Normal} \left( \log \tilde{R}_y, \phi_y^2 \right)$$

which leads to the definition of  $\ell_{SR}$  given above. In all cases  $a$  and  $b$  are strictly positive, and with the quantities  $F$ ,  $R$ , etc. linear models are used to parameterise  $\log a$  and/or  $\log b$ , where relevant.

By default, recruitment  $R$  as apposed to the recruitment predicted from a stock recruitment model  $\tilde{R}$ , is specified as a linear model with a parameter for each year, i.e.

$$\log R_y = \gamma_y$$

This is to allow modelled recruitment  $R_y$  to be shrunk towards the stock recruitment model. However, if it is considered appropriate that recruitment can be determined exactly by a relationship with covariates, it is possible, to instead define  $\log R$  in terms of a linear model in the same way as  $\log F$ ,  $\log Q$ ,  $\log \sigma$  and  $\log \tau$ .

## 4.1 Model fitting

The parameters that are estimated are log recruitment  $r_t$ , survey catchability  $q_a$ , the  $F$  parameters (MORE ON THESE LATER) and log stock numbers  $n_{a1}$  in the first year. The model is written in terms of these parameters,

$$c_{at} = r_{t-a+1} - \sum_{i=1}^{a-1} Z_{a-i, t-i} + \log \left( \frac{F_{at}}{Z_{at}} \left( 1 - e^{-Z_{at}} \right) \right) + \epsilon_{at}$$

$$s_{at} = q_a + r_{t-a+1} - \sum_{i=1}^{a-1} Z_{a-i, t-i} + \epsilon'_{at}$$

where  $\epsilon$  denotes the Gaussian observation error. There are some modifications required for the early cohorts as they use  $n_{a1}$  rather than recruits and for the plus groups, but these are trivial and not presented. It is possible to write these equations in matrix notation if we combine all catches into a single vector:  $\mathbf{c} = (c_{11}, c_{21}, \dots, c_{A+1,1}, c_{12}, \dots, c_{A+2,2}, \dots)^T$  and do similarly for the survey indices  $\mathbf{s}$ , it is also simpler if we define  $\mathbf{r} = (n_{A1}, n_{A-1,1}, \dots, n_{21}, r_1, \dots)^T$  and combine the  $F$  model parameters into a single vector  $\mathbf{f}$  then we can write the model as

$$\mathbf{c} = \mathbf{X}_r \mathbf{r} + o_1(\mathbf{f}) + o_2(\mathbf{f}) + \boldsymbol{\epsilon}$$

$$\mathbf{s} = \mathbf{M} \mathbf{X}_q \mathbf{q} + \mathbf{M} \mathbf{X}_r \mathbf{r} + \mathbf{M} o_1(\mathbf{f}) + \boldsymbol{\epsilon}'$$

where the functions  $o_1$  and  $o_2$  are nonlinear (vector) functions of the  $F$  model parameters and the  $\mathbf{X}$  matrices are various design matrices based on the full set of ages and years (to be described later). The  $\mathbf{M}$  matrix maps the correct age and year in the survey to that in the full set of ages and years. These equations can be further combined by stacking the equations

$$\begin{pmatrix} \mathbf{c} \\ \mathbf{s} \end{pmatrix} = \begin{pmatrix} \mathbf{0} & \mathbf{X}_r \\ \mathbf{M} \mathbf{X}_q & \mathbf{M} \mathbf{X}_r \end{pmatrix} \begin{pmatrix} \mathbf{q} \\ \mathbf{r} \end{pmatrix} + \begin{pmatrix} \mathbf{I} & \mathbf{I} \\ \mathbf{M} & \mathbf{0} \end{pmatrix} \begin{pmatrix} o_1(\mathbf{f}) \\ o_2(\mathbf{f}) \end{pmatrix} + \begin{pmatrix} \boldsymbol{\epsilon} \\ \boldsymbol{\epsilon}' \end{pmatrix}$$

so that the model is of the form

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + o(\mathbf{f}) + \boldsymbol{\epsilon}$$

where

$$\boldsymbol{\epsilon} \sim N \left( \mathbf{0}, \mathbf{W}^{-1} \right) \quad \text{where} \quad \mathbf{W}^{-1} = \begin{pmatrix} \sigma_c^2 \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \sigma_s^2 \mathbf{I} \end{pmatrix} = \sigma_c^2 \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \frac{\sigma_s^2}{\sigma_c^2} \mathbf{I} \end{pmatrix}$$

In other words, this statistical catch at age model can be written as a linear model with an offset due to nonlinear functions of the F model parameters. We use this to reduce the parameters in the fitting process by concentrating the likelihood. This can be done by inserting the maximum likelihood estimates of  $\beta$  conditional on  $\mathbf{f}$  and  $\mathbf{W}$  into the likelihood. The maximum likelihood estimate of  $\beta$  conditional on the other parameters is

$$\hat{\beta} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{y} - o(\mathbf{f}))$$

If we decide that the surveys and catches have the same observation variance then we can also estimate the observation variance conditionally, in this case

$$\begin{aligned} \hat{\beta}(\mathbf{f}) &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - o(\mathbf{f})) \\ \hat{\mathbf{y}}(\mathbf{f}) &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - o(\mathbf{f})) + o(\mathbf{f}) \\ &= \mathbf{H} (\mathbf{y} - o(\mathbf{f})) + o(\mathbf{f}) \\ &= \mathbf{H} \mathbf{y} + (1 - \mathbf{H}) o(\mathbf{f}) \end{aligned}$$

and the conditional (unbiased) estimate of  $\sigma^2$  is

$$\hat{\sigma}^2(\mathbf{f}) = \frac{1}{n - p} \log \left( (\hat{\mathbf{y}}(\mathbf{f}) - \mathbf{y})^T (\hat{\mathbf{y}}(\mathbf{f}) - \mathbf{y}) \right)$$

where  $n$  is the length of  $\mathbf{y}$  and  $p$  is the number of unique parameters in  $\beta$ . The concentrated log likelihood for  $\mathbf{f}$  is then

$$\begin{aligned} l(\mathbf{f}) &\propto -\frac{n}{2} \log \left( (\hat{\mathbf{y}}(\mathbf{f}) - \mathbf{y})^T (\hat{\mathbf{y}}(\mathbf{f}) - \mathbf{y}) \right) \\ &= -\frac{n}{2} \log \left( \left\| (1 - \mathbf{H}) \mathbf{y} + (1 - \mathbf{H}) o(\mathbf{f}) \right\|^2 \right) \end{aligned}$$

Since  $1 - \mathbf{H}$  and  $(1 - \mathbf{H}) \mathbf{y}$  only depend on the data these can be calculated outside of an iterative optimization procedure. It is straightforward to give different weights to the survey and catch components without increased computation (this is the same as assuming you know the ratio  $\frac{\sigma_c^2}{\sigma_s^2}$ ). However, if you want to estimate both variances then the estimate of  $\beta$  is

$$\begin{aligned} \hat{\beta}(\mathbf{f}) &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{y} - o(\mathbf{f})) \\ \hat{\mathbf{y}}(\mathbf{f}) &= \mathbf{X} (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{y} - o(\mathbf{f})) + o(\mathbf{f}) \end{aligned}$$

and as  $\mathbf{W}$  is not known before hand the inverse  $(\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1}$  must be computed at every iteration. The concentrated log likelihood in this case is

$$\begin{aligned} l(\mathbf{f}) &= \frac{1}{2} \log |\mathbf{W}| - \frac{1}{2} \left( (\hat{\mathbf{y}}(\mathbf{f}) - \mathbf{y})^T \mathbf{W} (\hat{\mathbf{y}}(\mathbf{f}) - \mathbf{y}) \right) \\ &= -\frac{n_c}{2} \log \sigma_c^2 - \frac{n_s}{2} \log \sigma_s^2 - \frac{1}{2} \left( (\hat{\mathbf{y}}(\mathbf{f}) - \mathbf{y})^T \mathbf{W} (\hat{\mathbf{y}}(\mathbf{f}) - \mathbf{y}) \right) \end{aligned}$$

## 5 Extending the model (brief)

First potential models for  $q_a$ . Linear forms can be included by simply changing the design matrices. By linear forms I mean spline smoothers with fixed degrees of freedom. An interesting addition would be to use penalised splines or better (I think) structured random effects

The model for F has not been considered so far. Options for this are a simple separable model, a model with several separable periods all these models can be expressed as linear models on the log link. Within the same framework is simple then to include smoothers (splines) with fixed degrees of freedom. More interesting but perhaps out with the scope of this project are structured random effect models for F. These include seasonal models (treating the number of ages as the season length) and correlated random walks.

## 6 Summary

At its simplest the model is a non-linear fixed effects regression, fitting the F parameters to the data but this requires that the survey variances are known relative to the catch variance. This model can include splines, or random effects with known parameters (variances, degrees of freedom, autocorrelation) in the q and r models. The reason this is being considered is that computational

speed is important and allowing some fixed variability may reduce bias in model outputs by allowing some flexibility, it is acknowledged that such assumptions make statistical testing a bit dodgy.

The next level of complexity is where we want to estimate the survey variance. This means we have to recalculate the hat matrix at every iteration which involves a matrix inversion. Since the matrix inversion is being done already, including structured random effects for  $q$  and recruitment into this only adds the variance parameters to the objective function.

If a stock recruit function were to be added, recruits could not be integrated out and would have to be estimated in the objective function.

If more complicated  $F$  models were used, such as structured random effects (random walks, correlated random walks, seasonal models) the number of parameters would increase in the objective function. Therefore a model with few parameters is one with a highly parametrised  $F$  model and no SRR relationship.

## 7 Implementation (sketch)

We require two functions that return an objective

A inputs are  $F$  at age and observation error and arguments are the data and the hat matrix

B inputs are  $F$  at age and any variance parameters taking as arguments the design matrix  $H$ , the weight matrix  $W$  and the structural prior matrix (not mentioned yet but lets call it  $Q$ )

The full objective function is then

1. take input parameters ( $F$  pars, variances, recruitments (if SRR model being used))
2. convert  $F$  pars into  $F$  at age
3. calculate objective value using one of the two functions A or B above
4. add on SRR density and prior densities for variances if necessary