

Maturity ogive for the southern hake stock

Cousido-Rocha, M., Izquierdo, F., Martinez-Minaya, J., Mendes, H., Silva, C., Silva, A.V., Cerviño, S.

2022-03-09

Contents

Model (theoretical explanation)	1
Exploratory	2
Motivation	6
Prepare data	8
Model total	9
Model by year	13
Supplementary material	17

Below the objective and a theoretical explanation of the model are reported. However the details of both can be find through the document tabs which explain the analysis step by step.

Objective: A combined maturity ogive (maturity proportions-at-length) for the southern hake stock estimated through the data derived from both institutes (laboratories), IPMA (Instituto Português do Mar e da Atmosfera) and IEO (Instituto Español de Oceanografía).

Model (theoretical explanation)

Maturity proportions-at-length have been estimated by bayesian regression models using the integrated nested Laplace approximation (INLA) (Rue et al., 2009) approach in the R-INLA software (<https://www.r-inla.org/>).

For estimating a combined maturity ogive for both laboratories a bivariate model has been required (Zuur and Ieno, 2018, additional details in Paradinas et al., 2017 and Izquierdo et al., 2021). The bivariate response variable is defined as follows.

$y_i^{IEO} \sim \text{Bernoulli}(\pi_i^{IEO})$, $i = 1, \dots, N^{IEO}$; being N^{IEO} the number of individuals measured by IEO.
 $y_j^{IPMA} \sim \text{Bernoulli}(\pi_j^{IPMA})$, $j = 1, \dots, N^{IPMA}$; being N^{IPMA} the number of individuals measured by IPMA.

The covariables (explanatory variables) are the length and the year. The length variable is introduced linear. On the other hand, the year covariable is introduced differently depending on the aim: a standard year combined maturity ogive (Approach 1) or a combined maturity ogive by year (Approach 2).

Approach 1

The year variability is taken into account through the random effect $a_i, a_j \sim N(0, \sigma_{year}^2)$, $i = 1, \dots, N^{IEO}$, $j = 1, \dots, N^{IPMA}$. Note that σ_{year}^2 parameter is common for IEO and IPMA response variables.

$$\text{Logit}(\pi_i^{IEO}) = \ln(\pi_i^{IEO}/(1 - \pi_i^{IEO})) = \beta_0 + \beta_1 \times (l^{IEO}(i)) + a_i + \epsilon_i$$

$$\text{Logit}(\pi_j^{IPMA}) = \ln(\pi_j^{IPMA}/(1 - \pi_j^{IPMA})) = \beta_0 + \beta_1 \times (l^{IPMA}(j)) + a_j + \epsilon_i$$

$l^{IEO}(i)$ assigns to each individual of IEO its corresponding length. The same for $l^{IPMA}(j)$. $\epsilon_i, \epsilon_j \sim N(0, \sigma_\epsilon^2)$; $a_i, a_j \sim N(0, \sigma_{year}^2)$.

Approach 2

The year is included in the model as a factor covariable.

$$\begin{aligned} \text{Logit}(\pi_i^{IEO}) &= \ln(\pi_i^{IEO}/(1 - \pi_i^{IEO})) = \beta_0 + \beta_1 \times (l^{IEO}(i)) + year_i + \epsilon_i \\ \text{Logit}(\pi_j^{IPMA}) &= \ln(\pi_j^{IPMA}/(1 - \pi_j^{IPMA})) = \beta_0 + \beta_1 \times (l^{IPMA}(j)) + year_j + \epsilon_i \end{aligned}$$

$l^{IEO}(i)$ assigns to each individual of IEO its corresponding length. The same for $l^{IPMA}(j)$. $year_i, year_j$ is a categorical covariate allowing for a different mean value per year. $\epsilon_i, \epsilon_j \sim N(0, \sigma_\epsilon^2)$.

References

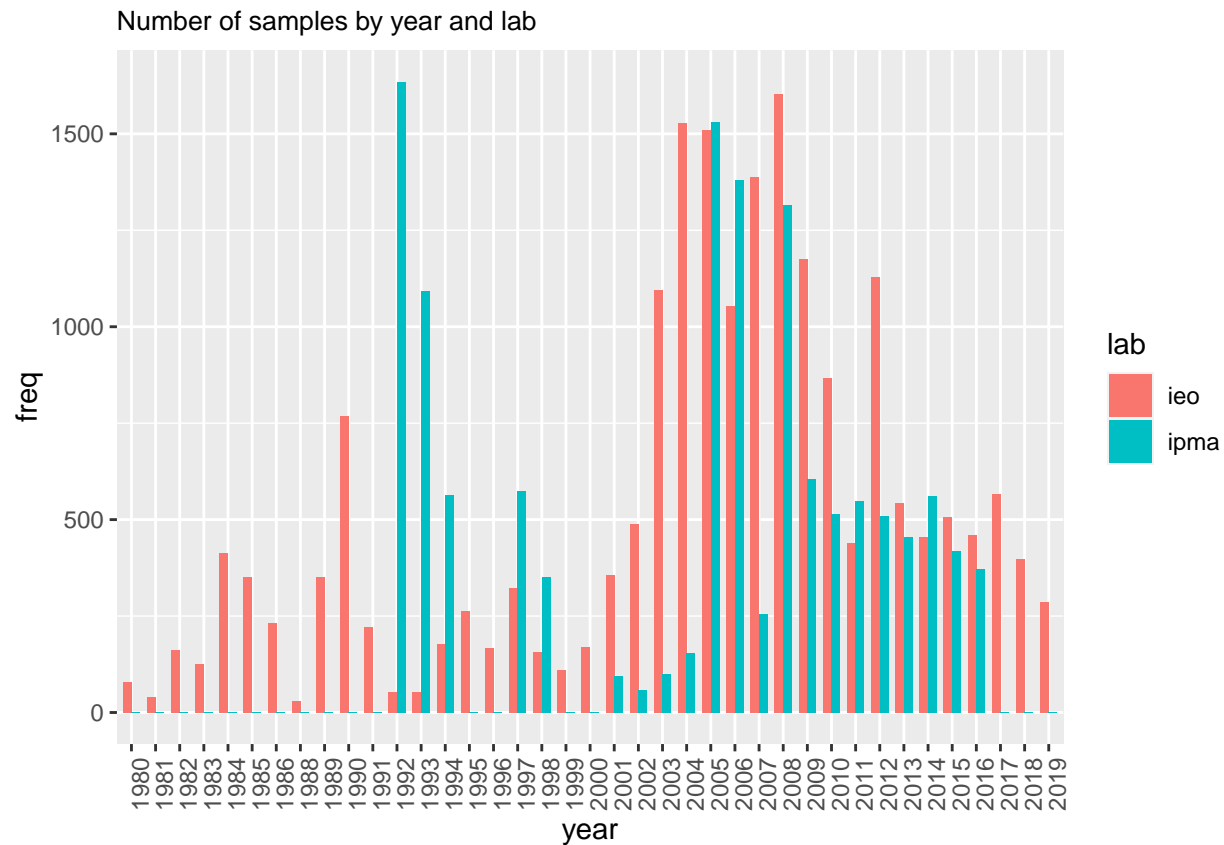
- Izquierdo, F., Paradinas, I., Cerviño, S., Conesa, D., Alonso-Fernández, A., Velasco, F., ... & Pennino, M. G. (2021). Spatio-temporal assessment of the European hake (*Merluccius merluccius*) recruits in the northern Iberian Peninsula. *Frontiers in Marine Science*, 8, 1.
- Paradinas, I., Conesa, D., Lopez-Quilez, A., & Bellido, J. M. (2017). Spatio-temporal model structures with shared components for semi-continuous species distribution modelling. *Spatial Statistics*, 22, 434-450.
- Rue, H., Martino, S., and Chopin, N. (2009). Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *J. R. Stat. Soc. B.* 71, 319–392. doi: 10.1111/j.1467-9868.2008.00700.x
- Wood, S.N. (2017) *Generalized Additive Models: An Introduction with R* (2nd edition). Chapman and Hall/CRC
- Zuur, A. F., Ieno, E. I. (2018). *Beginner's Guide to Spatial, Temporal and Spatial-Temporal Ecological Data Analysis with R-INLA Volume II: GAM and zero-inflated models* Published by Highland Statistics Ltd. Highland Statistics Ltd. Newburgh United Kingdom

Exploratory

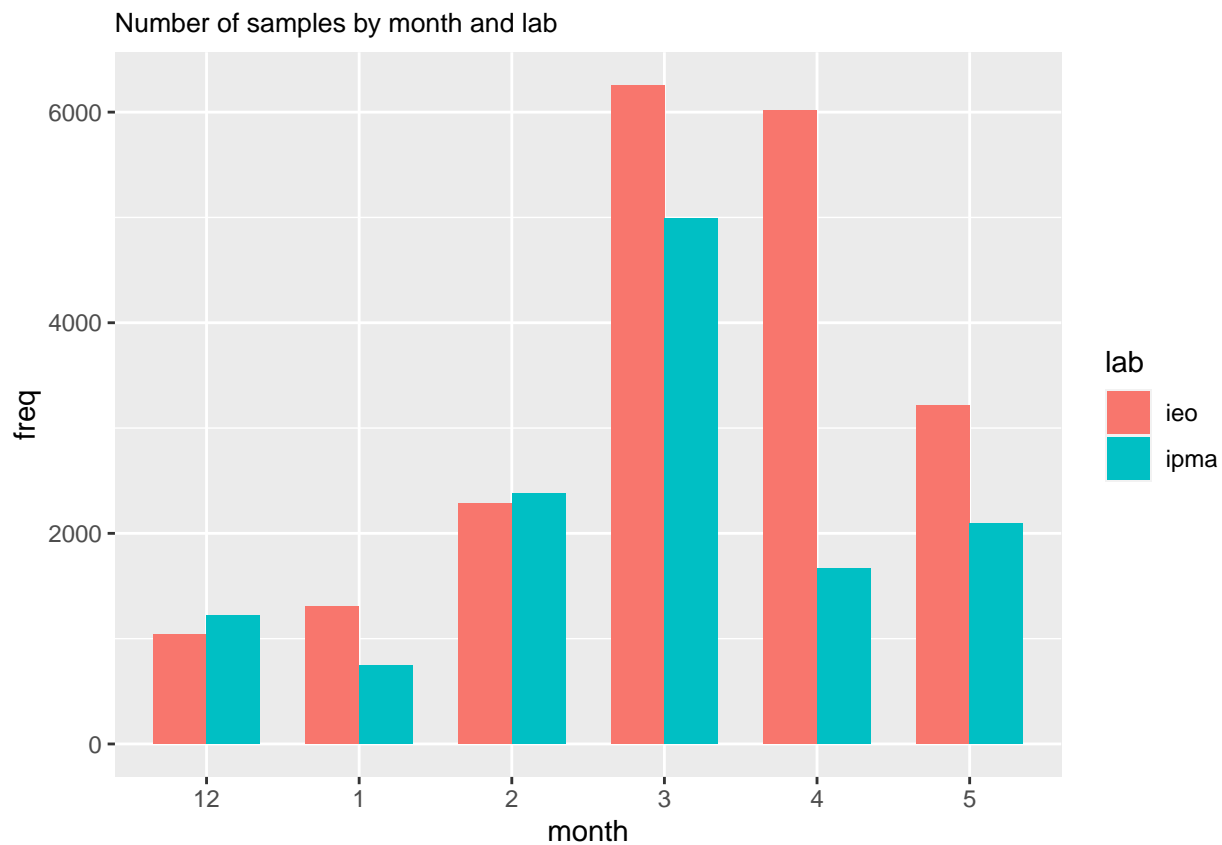
The data set contains the year of maturity, the month, the length (lt), the sex, the year of sample and the laboratory (institute) as you can see below. Note that for this study we have considered a subset of the data considering only females (sex=2).

```
##   year_mat month lt sex mat year_sample lab
## 1    1980     5 56  2   1    1980 ieo
## 2    1980     5 51  2   1    1980 ieo
## 3    1980     5 53  2   0    1980 ieo
## 4    1980     5 51  2   0    1980 ieo
## 5    1980     5 49  2   0    1980 ieo
## 6    1980     5 55  2   0    1980 ieo
```

The following plot report the number of samples for each year and institute. IPMA has no maturity data for the following years: 1980-1991, 1995, 1996, 1999, 2000, 2017-2019. IEO data is provided for the completed time period 1980-2019. Note that 2020 maturity data was provided only in May by the IEO. Since the information for this year is incomplete and may cause bias in the estimation of the ogive it has been decided to eliminate it.

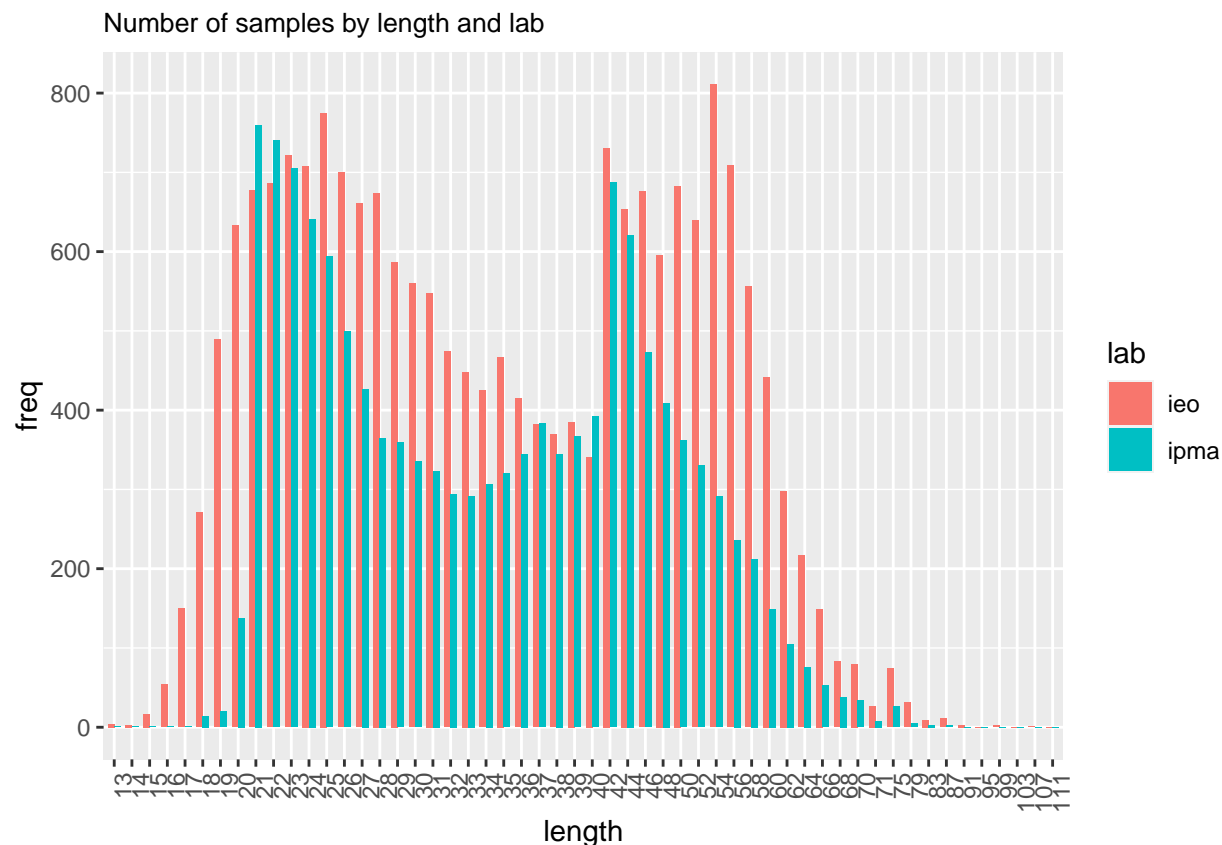


Next plot reports the number of samples by month and institute. Maturity data was compiled from the IEO and IPMA samples only for the spawning season, December to May. Note that, samples collected in December were allocated to the following year. Larger IPMA sampling corresponds to February and March, whereas for the IEO the larger sampling corresponds to March and April.



Next plot reports the number of samples by length and institute (laboratory). Overall good sampling of relevant length classes (from 20cm to 70cm).

```
##
##      13  14  15  16  17  18  19  20  21  22  23  24  25  26  27  28  29  30
## ieo    4   3  17  54 150 271 489 633 677 686 721 708 775 700 661 674 587 560
## ipma    1   1   1   1   1  14  20 138 760 740 705 641 594 499 427 365 359 335
##
##      31  32  33  34  35  36  37  38  39  40  42  44  46  48  50  52  54  56
## ieo  548 474 448 425 467 415 382 369 385 341 731 654 676 596 683 639 811 709
## ipma 323 294 292 307 320 345 383 345 367 392 687 621 473 409 362 330 291 236
##
##      58  60  62  64  66  68  70  71  75  79  83  87  91  95  99 103 107 111
## ieo  556 442 298 217 149  83  80  27  74  31   9  11   3   0   2   0   1   0
## ipma  212 149 105  76  53  38  34   8  26   5   3   2   0   0   0   0   0   0
```



Following 2010 benchmark it was decided to cut the ogive assigning zero to lengths below 21 cm because they are not mature.

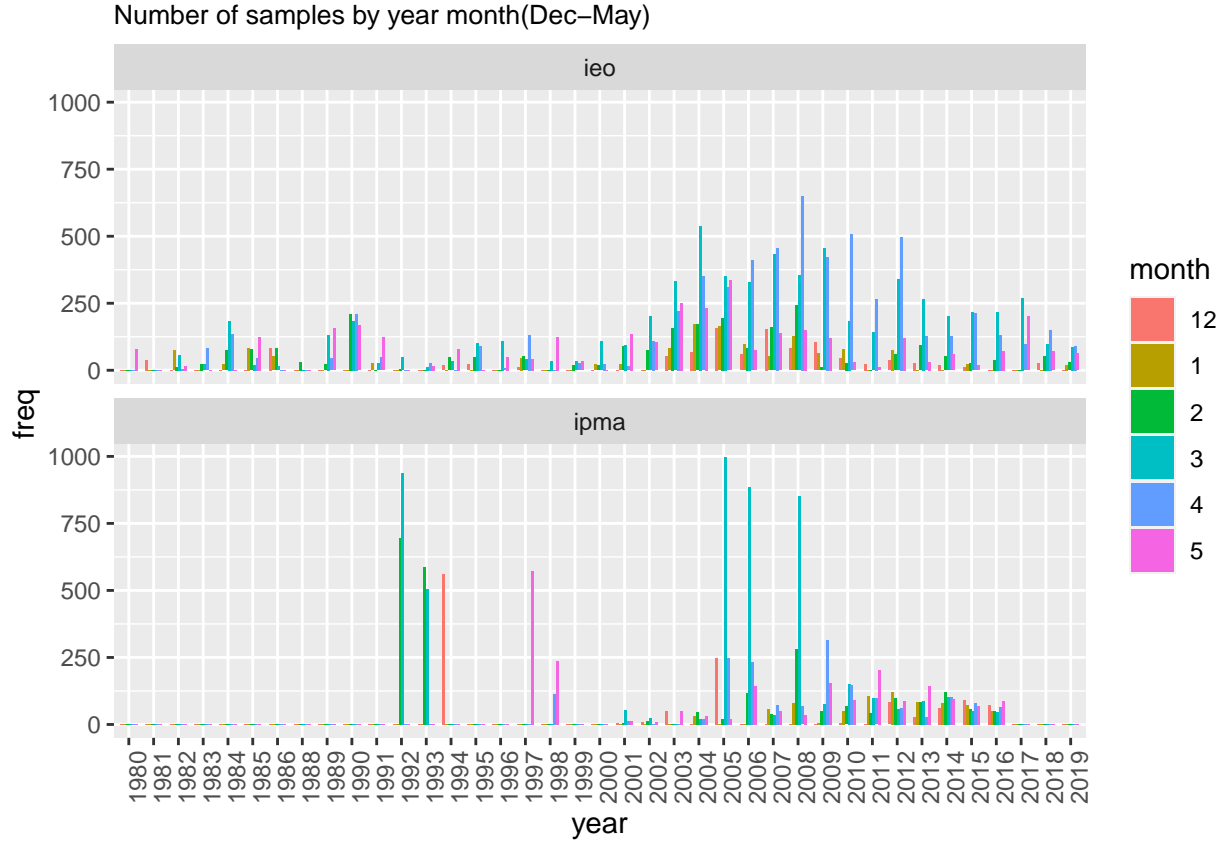
```
## [1] 15
```

```
##          lt mat
## 4336  17.5   1
## 8529  20.5   1
## 21020 18.7   1
## 21021 19.2   1
## 21024 19.6   1
## 21026 19.9   1
## 21027 20.7   1
## 21029 20.5   1
## 28788 19.6   1
## 28789 20.2   1
## 30068 19.1   1
## 34919 15.6   1
## 45245 17.7   1
## 51715 20.3   1
## 51873 18.2   1
```

Next plot reports the number of samples by year, month and institute. The plot shows that previously to 2001 IPMA information is missing except for 1992, 1993, 1994, 1997 and 1998. Furthermore, IEO sample size before 2001 is low and for some years not all months of the spawning season has been sampled. According to that years 1980-2000 are grouped for the modeling. On the other hand, for years 2017-2019 there are

not IPMA information and the IEO samples sizes are again low. Hence, such years are also grouped in the modeling.

Hence, our year covariable is not the year specific level factor is a year specific category factor with the following categories: 1980-2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017-2019.



Motivation

The maturity data is provided by two countries, Portugal and Spain, and a combined maturity ogive is required. Previous analysis provides evidences that in Portugal the maturity occurs at lower lengths than in Spain. In fact the regression logistic model (generalized linear model) below explains the maturity (binary response, immature/mature) using the length and the country factor leading to two statistical different ogives for each country.

The maturity data covers from 1980 to 2019, however, while the Spanish data cover the entire period, we have missing Portugal data for some years, and furthermore the samples sizes by year for each country are not balanced. For that reason the unification of the maturity data on an unique sample ignoring the country for further modeling, using for example glm, is not a suitable option. Other option can be a weighted average of the country ogives, but for that it is necessary to decide which weights must be used. After some research, we have found a possible solution using a Bayesian approach.

Our proposal is a bivariate bayesian regression model using the integrated nested Laplace approximation (INLA) (Rue et al., 2009) approach in the R-INLA software (<https://www.r-inla.org/>).

```

df2 <- data
mod.lab2 <- glm(mat ~ lt*lab, family = binomial(logit), data = df2)
summary(mod.lab2)

##
## Call:
## glm(formula = mat ~ lt * lab, family = binomial(logit), data = df2)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3.9076  -0.2908  -0.1078   0.1922   3.5054
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -12.061738   0.179328 -67.261  < 2e-16 ***
## lt           0.276627   0.004146  66.714  < 2e-16 ***
## labipma      1.482101   0.260561   5.688 1.28e-08 ***
## lt:labipma   -0.022793   0.006250  -3.647 0.000266 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 41136  on 33197  degrees of freedom
## Residual deviance: 15888  on 33194  degrees of freedom
## AIC: 15896
##
## Number of Fisher Scoring iterations: 7

```

#L50 Females IEO

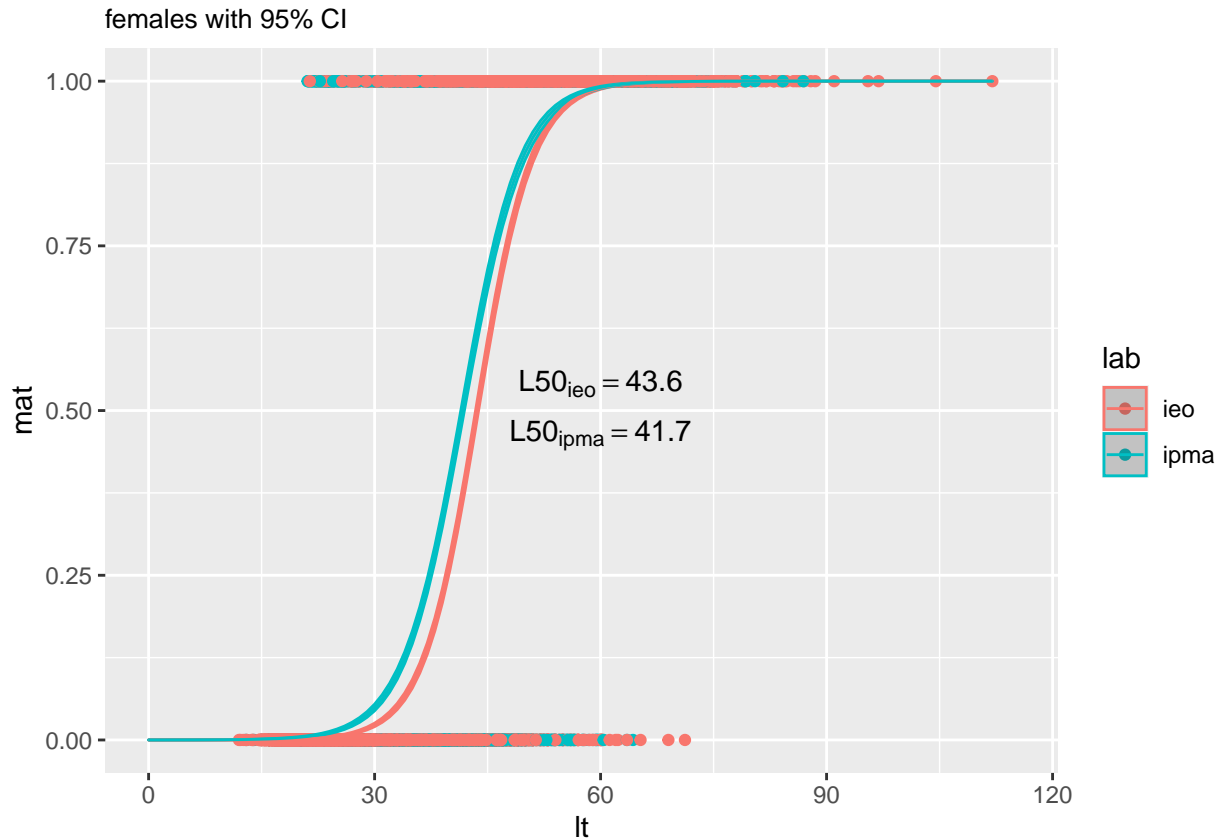
```
-(coef(mod.lab2)[1]/coef(mod.lab2)[2])
```

```
## (Intercept)
##      43.60289
```

#L50 Females IPMA

```
-(coef(mod.lab2)[1]+coef(mod.lab2)[3])/(coef(mod.lab2)[2]+coef(mod.lab2)[4])
```

```
## (Intercept)
##      41.67934
```



Prepare data

The bivariate model response considers separately two maturity variables one for each country. The two response variables are explained using length and year covariables. The model formulation in terms of covariables depends on the aim: - (i) a standard year combined maturity ogive or - (ii) a combined maturity ogive by year.

On (i) the common predictor for the two responses is equal to an intercept plus a linear effect of the length plus a year random effect. The year random effect is changed by a year factor for (ii) approach. The model carried out a combined estimation of all the parameters of the common predictor providing a combined maturity to introduce in the stock assessment model.

NOTE: as mentioned previously year covariable has the following categories: 1980-2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017-2019.

```
# Prepare data -----

NLbins<-c(seq(from=20, to=40, by=1),seq(from=42, to=70, by=2)) # Desired bins (SS model) 67
l_b=length(NLbins)

len=data$lt
l_len=length(len);aux=rep(0,l_len)

years<-(min(as.numeric(as.character(data$year_mat))):max(as.numeric(as.character(data$year_mat))))
```



```

# Response -----

data_ieo=subset(data,data$lab=="ieo")
data_ipma=subset(data,data$lab=="ipma")
data=rbind(data_ieo,data_ipma)

ind_ieo=which(data$lab=="ieo")
ind_ipma=which(data$lab=="ipma")
len=length(data$lab)

len_ieo=length(ind_ieo)
len_ipma=length(ind_ipma)

YCombined <- matrix(NA, nrow = len, ncol = 2)
YCombined[1:len_ieo, 1] <- (data$mat[ind_ieo])
YCombined[(len_ieo+1):(len_ipma+len_ieo), 2] <- (data$mat[ind_ipma])

# Grouped years -----

# Years previous to 2001 into a group -----

data$Gyear_mat=as.character(data$year_mat)
ind=which(as.numeric(as.character(data$year_mat))<2001)
data$Gyear_mat[ind]="1980-2000"

# Years 2017,2018 and 2019 into a group -----
ind=which(as.numeric(as.character(data$year_mat))>2016)
data$Gyear_mat[ind]="2017-2019"

data$Gyear_mat=as.factor(data$Gyear_mat)

```

Model total

Standard ogive: a single ogive for both institutes and years.

Code

```

# Model 1 -----

f3 <- YCombined ~ 1 + lt +
      f(Gyear_mat, model = "iid")

I3 <- inla(f3,
  control.compute = list(config=TRUE,
                        dic = TRUE,
                        cpo=TRUE),
  family = c("binomial","binomial"),
  data = data,
  control.inla = list(strategy = 'adaptive'),
  verbose=TRUE, num.threads = 1)

```

```
summary(I3)
```

```
##
## Call:
##   c("inla(formula = f3, family = c(\"binomial\", \"binomial\"), data =
##     data, \"\", \" verbose = TRUE, control.compute = list(config = TRUE, dic =
##     TRUE, \"\", \" cpo = TRUE), control.inla = list(strategy = \"adaptive\"),
##     \"\", \" num.threads = 1)\")
## Time used:
##   Pre = 0.73, Running = 55.9, Post = 2.4, Total = 59.1
## Fixed effects:
##           mean      sd 0.025quant 0.5quant 0.975quant      mode kld
## (Intercept) -11.242 0.147    -11.534   -11.241    -10.955  -11.239    0
## lt           0.265 0.003      0.259    0.265      0.271    0.265    0
##
## Random effects:
##   Name      Model
##   Gyear_mat IID model
##
## Model hyperparameters:
##           mean      sd 0.025quant 0.5quant 0.975quant      mode
## Precision for Gyear_mat 13.66 5.08      5.90    12.93      25.61 11.54
##
## Expected number of effective parameters(stdev): 16.94(0.656)
## Number of equivalent replicates : 1959.42
##
## Deviance Information Criterion (DIC) .....: 15834.58
## Deviance Information Criterion (DIC, saturated) ....: 15834.57
## Effective number of parameters .....: 17.20
##
## Marginal log-Likelihood: -7947.11
## CPD and PIT are computed
##
## Posterior marginals for the linear predictor and
## the fitted values are computed
```

```
#INLAutils::plot_fixed_marginals(I3)
```

```
#INLAutils::plot_hyper_marginals(I3)
```

```
#INLAutils::plot_random_effects(I3)
```

```
# Prediction IPS -----
```

```
I1=I3
```

```
r=I3
```

```
r.samples = inla.posterior.sample(1000, r)
```

```
psam <- sapply(r.samples, function(x) {
```

```
  lt_effect <- x$latent %>% rownames(.) %>% stringr::str_detect("^lt") %>% x$latent[,]
```

```
  intercept <- x$latent %>% rownames(.) %>% stringr::str_detect("^\\(Intercept\\)") %>% x$latent[,]
```

```
  year_effect <- rnorm(length(lt_effect), sd = 1/sqrt(x$hyperpar[1]))
```

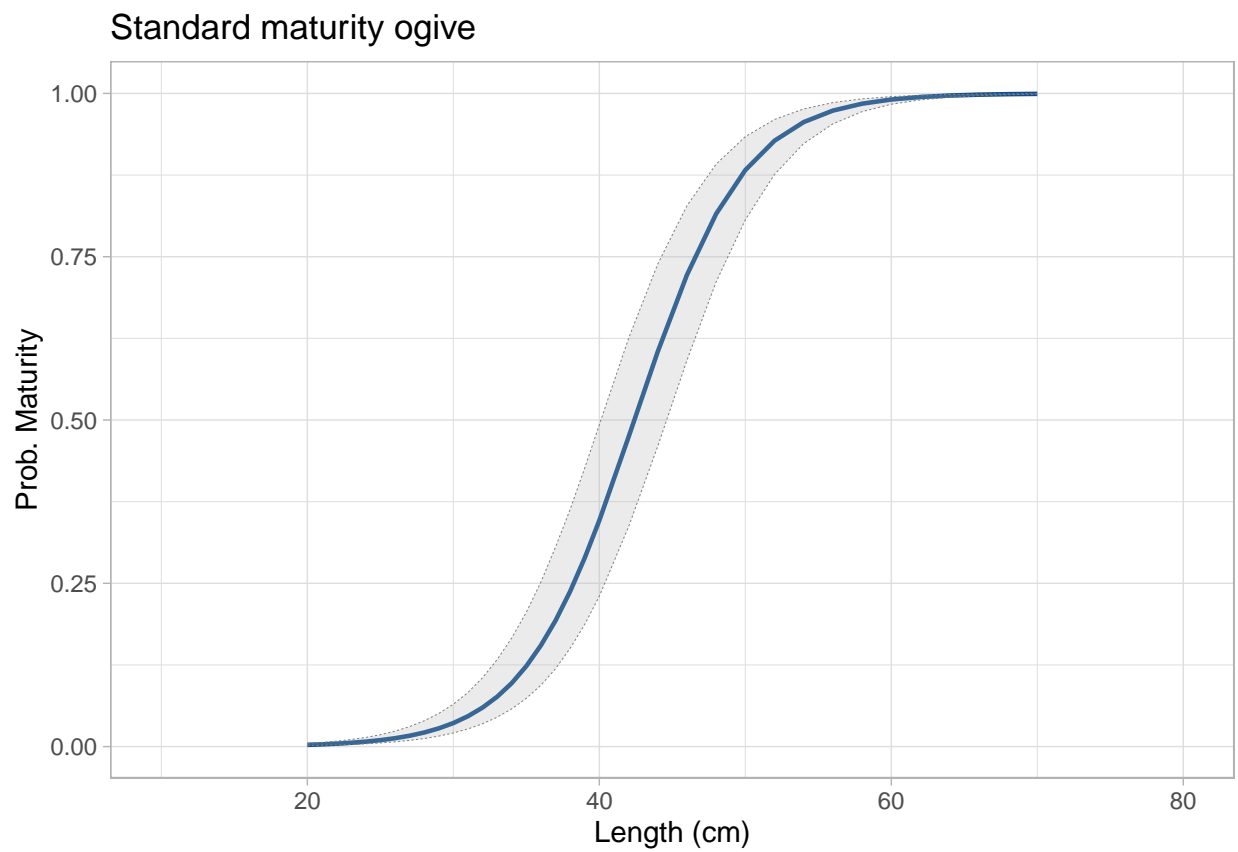
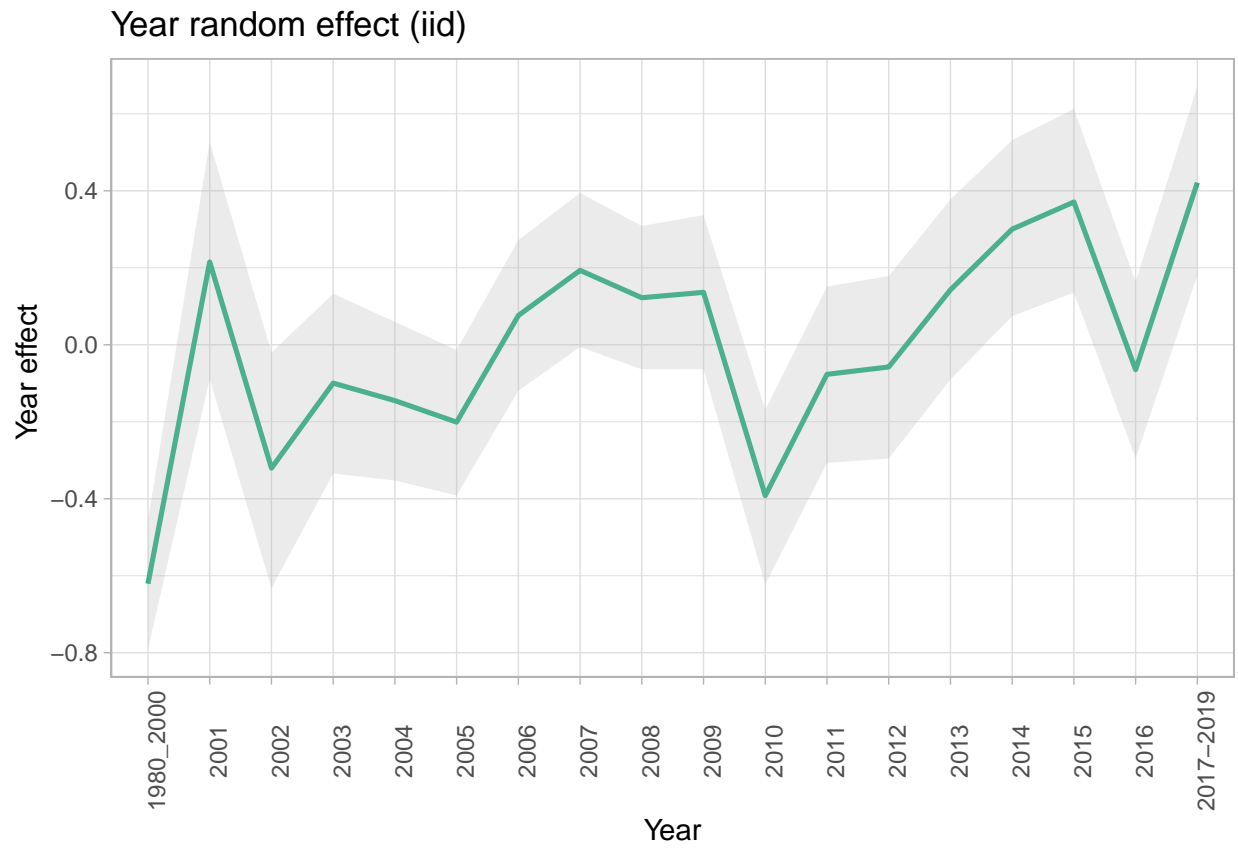
```
  predictor <- intercept + year_effect + lt_effect*NLbins
```

```
  exp(predictor)/(1 + exp(predictor))
```

```
})
```

```
q.sam_al_a <- apply(psam, 1, quantile,  
  c(.025, 0.05, 0.5, 0.95, .975), na.rm =TRUE)
```

Plot



L_{50} values

Length at 50% maturity.

```
##           L50      lower      upper
## 1 42.36566 40.35277 44.47734
```

Model by year

Yearly ogive: a specific ogive for year category.

Code

```
# Model 2 -----

f3 <- YCombined ~ 1 + lt + Gyear_mat

I3 <- inla(f3,
          control.compute = list(config=TRUE,
                                dic = TRUE,
                                cpo=TRUE),
          family = c("binomial","binomial"),
          data = data,
          control.inla = list(strategy = 'adaptive'),
          verbose=TRUE, num.threads = 1)
summary(I3)

##
## Call:
## inla(formula = f3, family = c("binomial", "binomial"), data =
## data, verbose = TRUE, control.compute = list(config = TRUE, dic =
## TRUE, cpo = TRUE), control.inla = list(strategy = "adaptive"),
## num.threads = 1)
## Time used:
## Pre = 0.373, Running = 18.1, Post = 2.91, Total = 21.4
## Fixed effects:
##              mean      sd 0.025quant 0.5quant 0.975quant      mode kld
## (Intercept) -11.912 0.142   -12.193  -11.911   -11.636 -11.909  0
## lt           0.266 0.003    0.260   0.266    0.272  0.266  0
## Gyear_mat2001 0.938 0.177    0.590   0.938    1.287  0.938  0
## Gyear_mat2002 0.203 0.172   -0.133   0.203    0.541  0.202  0
## Gyear_mat2003 0.523 0.117    0.295   0.523    0.752  0.522  0
## Gyear_mat2004 0.480 0.096    0.292   0.480    0.668  0.480  0
## Gyear_mat2005 0.427 0.083    0.263   0.427    0.590  0.427  0
## Gyear_mat2006 0.724 0.089    0.549   0.724    0.898  0.724  0
## Gyear_mat2007 0.850 0.092    0.670   0.850    1.030  0.850  0
## Gyear_mat2008 0.771 0.082    0.611   0.771    0.932  0.771  0
## Gyear_mat2009 0.789 0.092    0.608   0.789    0.969  0.789  0
## Gyear_mat2010 0.196 0.111   -0.022   0.196    0.412  0.196  0
## Gyear_mat2011 0.553 0.114    0.329   0.553    0.776  0.554  0
## Gyear_mat2012 0.574 0.120    0.338   0.574    0.809  0.574  0
```

```
## Gyear_mat2013      0.806 0.118      0.575  0.806      1.036  0.806  0
## Gyear_mat2014      0.985 0.113      0.763  0.985      1.205  0.985  0
## Gyear_mat2015      1.073 0.120      0.837  1.073      1.307  1.073  0
## Gyear_mat2016      0.567 0.114      0.342  0.568      0.791  0.568  0
## Gyear_mat2017-2019 1.137 0.124      0.893  1.137      1.380  1.137  0
##
## Expected number of effective parameters(stdev): 19.02(0.00)
## Number of equivalent replicates : 1745.93
##
## Deviance Information Criterion (DIC) .....: 15835.94
## Deviance Information Criterion (DIC, saturated) .....: 15835.93
## Effective number of parameters .....: 19.02
##
## Marginal log-Likelihood: -8008.40
## CPD and PIT are computed
##
## Posterior marginals for the linear predictor and
## the fitted values are computed
```

```
# Prediction IPS -----
I2=I3
r=I3
r.samples = inla.posterior.sample(1000, r)
psam <- sapply(r.samples, function(x) {

  lt_effect <- x$latent %>% rownames(.) %>% stringr::str_detect("^lt") %>% x$latent[,]
  intercept <- x$latent %>% rownames(.) %>% stringr::str_detect("^\\(Intercept\\)") %>% x$latent[,]
  beta_y <- x$latent %>% rownames(.) %>% stringr::str_detect("^Gyear_mat") %>% x$latent[,]

  predictor1990 <- intercept + lt_effect*NLbins

  pre=list();l=length(beta_y)
  for (i in 1:l){
    pre[[i]]=intercept + beta_y[i] + lt_effect*NLbins
  }

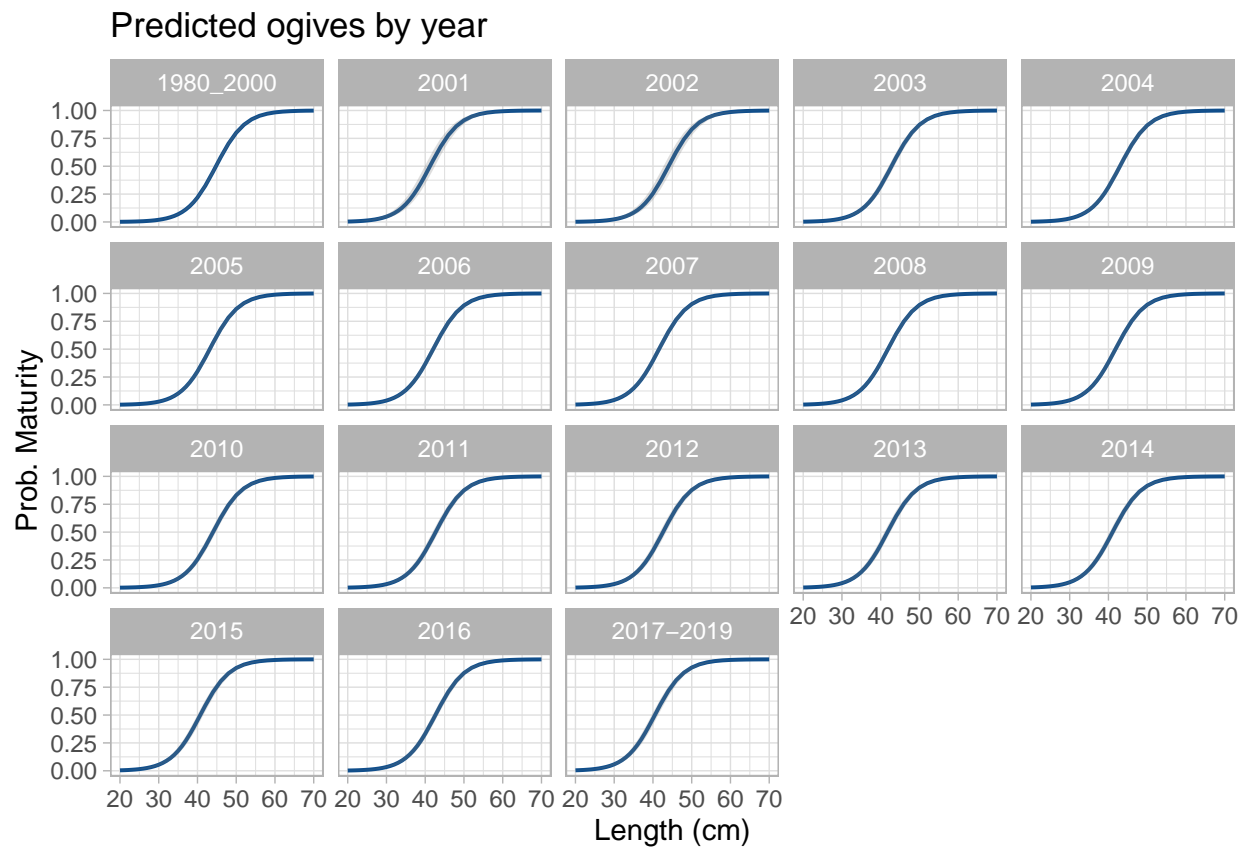
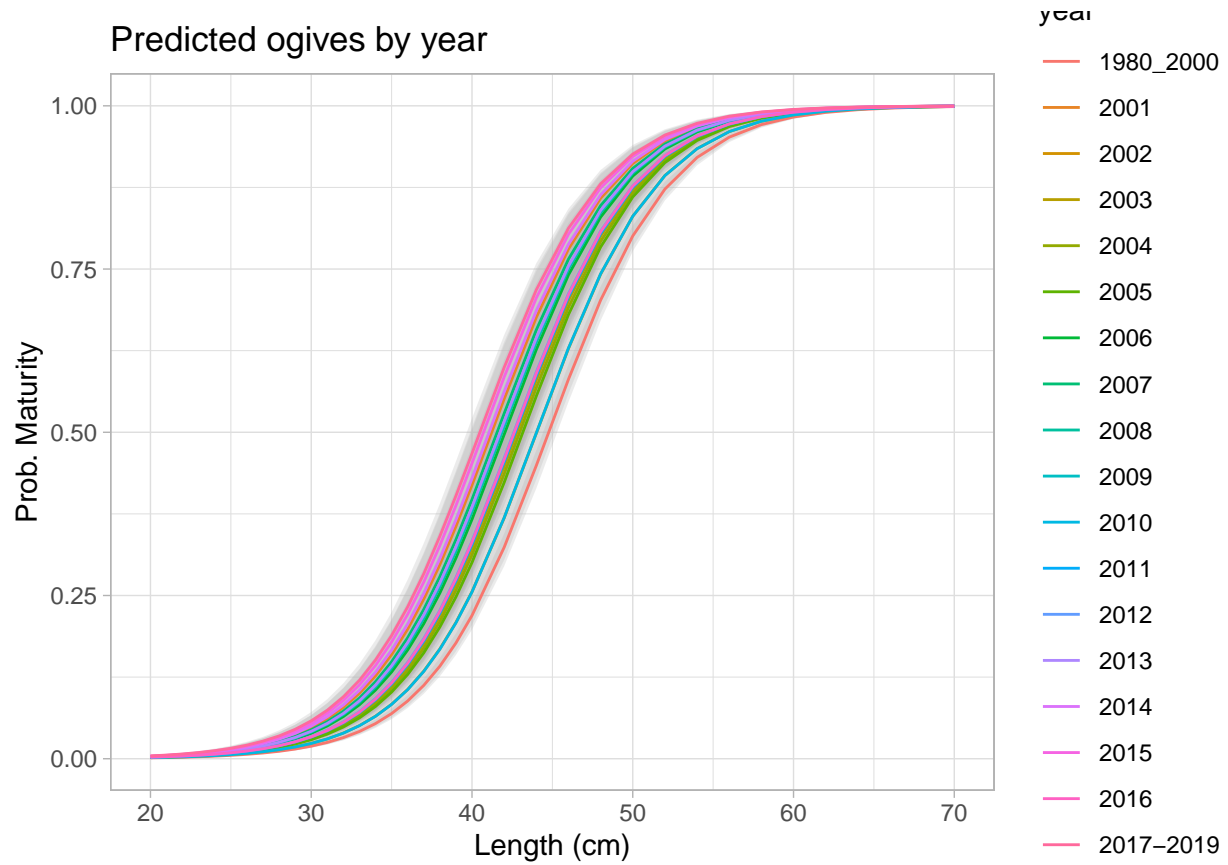
  predictor=predictor1990

  for (i in 1:l){
    predictor <- c(predictor, pre[[i]])
  }

  exp(predictor)/(1 + exp(predictor))
})

q.sam_al_a <- apply(psam, 1, quantile,
                    c(.025, 0.05, 0.5, 0.95, .975), na.rm =TRUE)
```

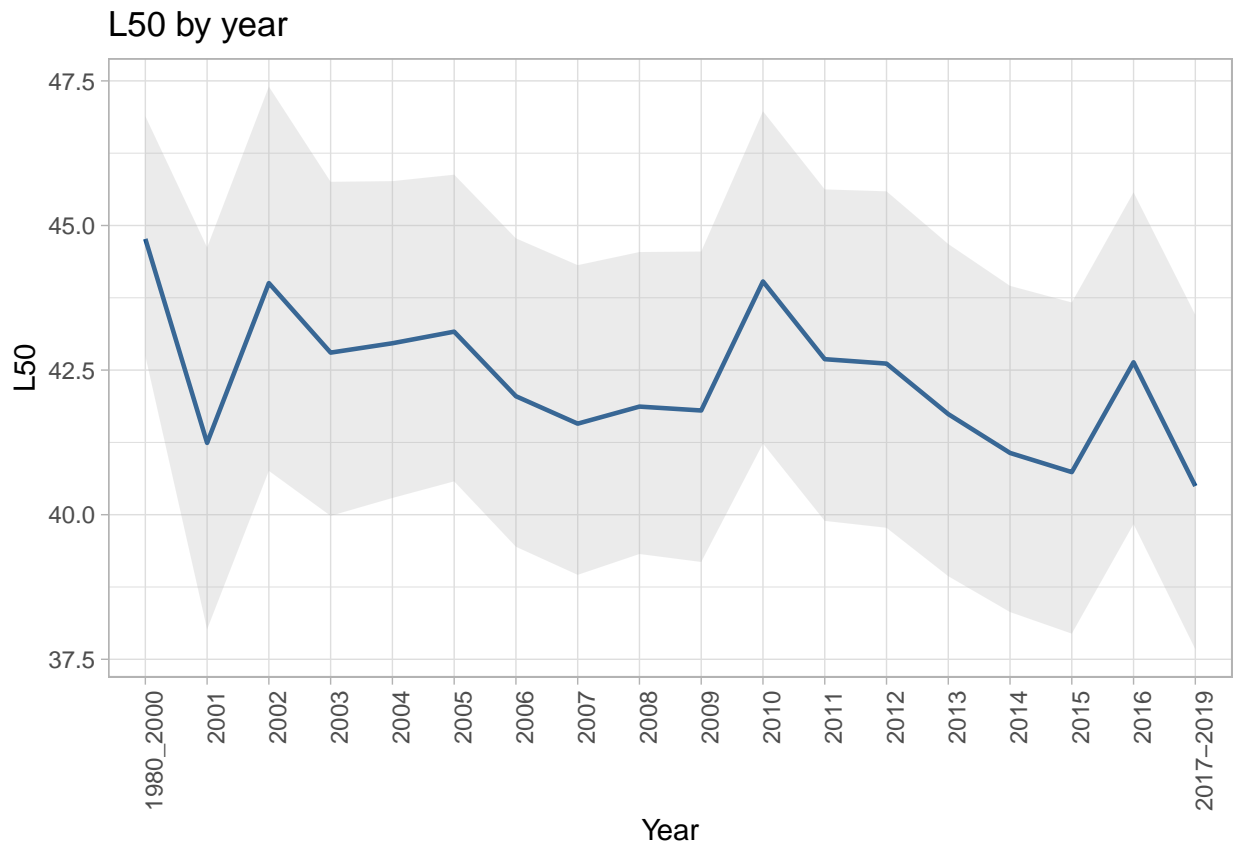
Plot



L_{50}

L_{50} (length at 50% maturity) times series. Since the analysis of the series shows clear variability among year categories, the time specific model is proposed to be used instead to the standard year combined maturity ogive.

##	L50	lower	upper	year
## 1	44.76778	42.74336	46.88983	1980_2000
## 2	41.24082	38.01732	44.62072	2001
## 3	44.00448	40.75780	47.40163	2002
## 4	42.80211	39.98009	45.75557	2003
## 5	42.96373	40.28891	45.76701	2004
## 6	43.16444	40.57634	45.87834	2005
## 7	42.04829	39.44575	44.77822	2006
## 8	41.57379	38.95971	44.31419	2007
## 9	41.86873	39.32035	44.54118	2008
## 10	41.80167	39.18236	44.54971	2009
## 11	44.03250	41.22972	46.97485	2010
## 12	42.68880	39.89288	45.62410	2011
## 13	42.61215	39.77339	45.59071	2012
## 14	41.73883	38.93614	44.67973	2013
## 15	41.06759	38.31661	43.95626	2014
## 16	40.73617	37.94327	43.66906	2015
## 17	42.63505	39.83871	45.57335	2016
## 18	40.49509	37.67432	43.45491	2017-2019



Supplementary material

Structural changes

A structural change analysis has been applied over the year time series of L_{50} (derived from the model using year factor covariable with a specific level for each year instead of the year categories). As you can see this analysis also reports 2000 as a break point of the time series in accordance with our conclusion after the exploratory analysis.

```
library(strucchange)

load("50.RData")

maturity<-dL50

# Input NA's (if is required)

interpFun <- function(dat) {
  for (i in 1:length(dat)){
    if (is.na(dat[i]))
      if(i == 1) {
        dat[i] <- rnorm(1,mean(dat, na.rm=T),
                        sd(dat, na.rm=T))
      } else {
        dat[i] <- rnorm(1,mean(dat[c(i-1, i+1)],na.rm=T),
                        sd(dat[c(i-1, i+1)],na.rm=T))
      }
  }
  return(dat)
}

# Define time series -----
mInterp <- interpFun(maturity$L50)
mInterp <- ts(mInterp,
              start=min(maturity$year),
              frequency = 1)

# Detect break and test-----
ocusm <- efp(mInterp~1, type="OLS-CUSUM")
#ocusm <- efp(mInterp~1, type="Rec-CUSUM")
#ocusm <- efp(mInterp~1, type="Rec-MOSUM")
#ocusm <- efp(mInterp~1, type="OLS-MOSUM")
bpm <- breakpoints(mInterp~1)
maturity$year[bpm$breakpoints]

## [1] 2000

sctest(ocusm)

##
## OLS-based CUSUM test
##
```

```
## data:  ocusm
## S0 = 1.4989, p-value = 0.02237
```

```
# Plot series + break point -----
```

```
plot(mInterp,
     xlab= "Year",
     ylab= "L50 parameter",
     lty=1,
     lwd=2,
     main = "L50 breakpoints")
lines(mInterp,
      lty = 1,
      lwd = 2)
abline(v=maturity$year[bpm$breakpoints],
      lwd= 2,
      lty = 1,
      col="blue")
legend("topright",
      legend = c("a", "bp L50"),
      lty = c(1,1),
      col = c("black", "blue"), bty="n", x.intersp=0.5, horiz= F, cex=0.70)
```

