# Modelling Spatio-temporal Abundance at Age with Bayesian Geostatistics and Compositional Data Analysis

Ernesto Jardim[1] <ernesto@ipimar.pt>

Paulo J. Ribeiro Jr[2] <paulojus@ufpr.br>

April 18, 2008

[1]Instituto Nacional de Investigação Agrária e das Pescas, Av. Brasilia, 1449-006, Lisboa, Portugal

[2]Laboratório de Estatística e Geoinformação, Universidade Federal do Paraná, C.P. 19.081 CEP: 81.531-990 , Curitiba, Paraná, Brasil

## Abstract

This work presents a methodology to estimate abundance at age by year combining the spatial distribution of the stock and the age structure in a single parametric model. By separating the age compositions from the age-aggregated abundance, suitable models can be applied to each variable improving the analysis of the data and increasing the flexibility of the model. The parametric characteristics of the model allows the usage of Monte Carlo methods, providing means to overcome difficulties in obtaining the analytical expression of abundance at age. On the other hand, Monte Carlo simulations can be used as inputs for large simulation frameworks like those use for Management Strategies Evaluations. Age structures were studied by compositional data analysis allowing the full covariance structure of age compositions to be considered. Age-aggregated observations were modelled with geostatistical methods explicitly modelling the correlation between abundance at different locations. The methodology produces abundance indicators that provide an overview of abundance along different perspectives. The analysis of age compositions provides an insight on how the population structure evolves over time. The geostatistical submodel returns abundance indicators for both, space and time dimensions. An application to Hake (*Merluccius merluccius*) caught by the Portuguese Bottom Trawl Surveys is presented, and methods are proposed to handle specific characteristics of the problem at hand. We suggest a calibration of the different conditions on which data were collected using a GLM with negative binomial distribution and several covariates which also deals with asymmetry and over-dispersion.

**Key-words:** abundance at age, bottom trawl survey, hake, geostatistics, compositional data analysis

# 1  Introduction

Estimates of abundance are important indicators of stock size and space-time distribution of marine populations. Such indicators contain valuable information for stock assessment, where they are used as fisheries-independent inputs, and, more generally, for fisheries advice and ecological management. Several methods have been proposed to study abundance using design-based techniques (Cochran 1960; Thompson 1992; Smith and Gavaris 1993); specific statistical distributions like log-normal (McConnaughey and Conquest 1993; Brynjarsdottir and Stefansson 2004; Dingsor 2005; Smith 1990), delta (Pennington 1983; Stefansson 1996; Smith 1988), Poisson and negative binomial (O'Neill and Faddy 2003; Pradhan and Leung 2006) or zero inflated distributions (Martin et al. 2005; Mendes 2007); and different modelling procedures like generalised linear models (Smith 1990; Stefansson 1996; Brynjarsdottir and Stefansson 2004; Chen et al. 2004; Sousa et al. 2007), generalised additive models (Piet 2002), geostatistics (Rivoirard et al. 2000; Roa-Ureta and Niklitschek in press) or hierarchical models (Mendes 2007).

Sampling fish populations will naturally originate data sets with high correlation, both in population structure and spatial distribution, once individuals with similar ages or lengths will assemble looking for the best geophysical conditions. Following the work on statistical analysis for compositional data by Aitchison (1982, 2003), Hrafnkelsson and Stefansson (2004) and Babak et al. (2007) describe methods to model the correlation between length groups using Bayesian methods and maximum likelihood estimators, respectively. Within this approach, population age structure is represented by compositional data, defined by vectors of proportions at age subject to the constraint of summing one. Spatial patterns encountered on abundance data are expressed by the correlation between observations related to the distance between the geographical locations where the observations were collected and modelled with geostatistical methods (Cressie 1993; Diggle et al. 1998; Chilès and Delfiner 1999; Diggle and Ribeiro 2007).

Our aim with this work is to propose a methodology combining the spatial distribution of the stock and the relation between age groups into a single model. The methodology provides a framework to obtain simulations of abundance at age that can be used as input to large simulation frameworks like Management Strategy Evaluation (MSE) (Hammond and Donovan in press; Johnston and Butterworth 2005; Punt et al. 2005; Kell et al. 2007), a major subject for modern scientific advice on fisheries and ecological management. An application to hake (*Merluccius merluccius*) caught by the Portuguese Bottom Trawl Survey (BTS) is presented, and methods to handle specific characteristics of modelling hake's abundance are proposed.

The next section describes the Portuguese BTS and the data set used for analysis. On the Methods section we will start by presenting the model and its most important characteristics followed by a detailed description of parameter estimation for abundance at age. The Results section describes the adjustments required to apply the proposed model to estimate hake's abundance at age and presents different perspectives of abundance:

the time series of age aggregated abundance showing the trends in biomass over time; the yearly spatial distribution of biomass showing areas of higher density of hake; and the yearly abundance at age which constitutes a major input parameter for stock assessment. Finally, we discuss the model and its limitations, and compare the results obtained with the abundance at age estimates obtained using design-based statistics.

## 2  Material

The Portuguese BTS have been carried out in Portuguese continental waters since 1979 on board the R/V Noruega and R/V Capricórnio. The main objectives of these surveys are: (i) estimate indices of abundance and biomass of the most important commercial species; (ii) describe the spatial distribution of the most important commercial species, and (iii) collect individual biological parameters such as maturity, sex-ratio, weight, food habits, etc. The target species are hake (*Merluccius merluccius*), horse mackerel (*Trachurus trachurus*), mackerel (*Scomber scombrus*), blue whiting (*Micromessistius poutassou*), megrims (*Lepidorhombus boscii* and *L. whiffiagonis*), monkfish (*Lophius budegassa* and *L. piscatorius*) and Norway lobster (*Nephrops norvegicus*). A Norwegian Campbell Trawl 1800/96 (NCT) with a codend of 20 mm mesh size, mean vertical opening of 4.8 m and mean horizontal opening between wings of 15.6 m has been used (Anonymous 2002).

A stratified sampling design was used to define locations for data collection between 1989 and 2004. The stratification was defined by 12 sectors along the Portuguese continental coast subdivided into 4 depth ranges: 20-100m, 101-200m, 201-500m and 501-750 m, with a total of 48 strata. Constraints in vessel time limited the sample size to 97 locations, evenly allocated to obtain two locations within each stratum. The coordinates of the sampling locations were selected randomly, albeit constrained by the historical records of clear tow positions and other information about the sea floor, avoiding places where trawling was not possible. In 2005 a new sampling design, composed by a regular grid with a set of additional random locations, was introduced following Jardim and Ribeiro Jr. (2007). The tow duration was 60 minutes until 2001 and then reduced to 30 minutes for the subsequent years, based on an experiment that showed no significant differences in the mean abundance and length distribution between the two tow durations (Cardador, pers.comm.). Historically the Portuguese Autumn bottom trawl survey has been carried out between September and December and hauls occurred during daylight. The number of hauls per year, the estimates of abundance by year together with its standard deviation and coefficient of variation are presented in the first five columns of Table 1. Sampling statistics of abundance at age per year and coefficient of variation are showed on the top panel of Table 2.

The data set included all valid hauls executed during the Autumn survey between 1989 and 2006. Each record corresponds to hake catches in number of individuals by age, haul duration (minutes), haul time, haul date, coordinates (UTM, Zone 29), bottom salinity and bottom temperature. Catches obtained with R/V Capricórnio (1996, 1999, 2003 and 2004) were calibrated to R/V Noruega's catches using factors by

3

85  age estimated in a calibration exercise in 2006 (Cardador, pers.comm). Figure 1 shows the map of observed

86  age aggregated catches of hake during the study period.

# 3    Methods

88  The main target of the analysis is to model the abundance at age, $I$, which is given by the product of two

89  random variables $I_{ij} = Y_i P_{ij}$ where $Y_i$ represents the age aggregated abundance for the $i^{th}$ year, $i = 1, \ldots, n$,

90  and $P_{ij}$ refers to the proportion of individuals at the $i^{th}$year and $j^{th}$age, $j = 1, \ldots, m$. The age composition

91  for each year is denoted by $\mathbf{P}_i$. The model aims to disentangle population abundance from the composition

92  by age, so that both quantities can be modelled independently and taking into account the nature of each

93  one. Monte Carlo methods combine outputs of both submodels to obtain samples of the distribution of $I$

94  allowing for inferences about $I_{ij}$. This section provide details on the models and methods adopted.

95  Observed data on abundance at age consists of the total catch per unit effort in year $i$, age $j$ and haul

96  $h = 1, \ldots, H$ represented by $C_{ijh}$, from which proportion at age is computed by $P_{ijh} = C_{ijh}(\sum_{j=1}^{m} C_{ijh})^{-1}$.

97  Compositional data analysis (Aitchison 1982, 2003) is used to model $\mathbf{P}_i$, transforming the $m$ proportions

98  $P_{ijh}$ to $m - 1$ additive log-ratios compositions $D_{ijh} = \log(P_{ijh} P_{ij=a,h}^{-1})$ with $j \neq a$. This is a convenient

99  scale for parameter estimation and simulation given that compositions follow approximately a multivariate

100 Gaussian distribution, $\mathbf{D}_i \sim MVG(\Lambda_i, \Sigma_i)$, from which the mean estimator estimator $\hat{\Lambda}_i \sim \text{MGV}(\mu_i, \varsigma_i)$.

101 The covariance structure of the age compositions can be estimated from the data and subsequently used

102 in the simulation procedure. Maximum likelihood estimators are given by $\hat{\mu}_i = \bar{\mu}_i$, the vector of marginal

103 arithmetic means, and $\hat{\varsigma}_i = \hat{\rho}(\mathbf{D}_i)\hat{\boldsymbol{\sigma}}_i^2 H_i^{-1}$, where $\hat{\rho}(\mathbf{D}_i)$ is the sample correlation matrix and $\hat{\boldsymbol{\sigma}}_i^2$ is the vector

104 of marginal sample variances (Murteira 1990). Parametric bootstrap (Efron and Tibshirani 1993) is used

105 to assess the variability of the proportions by sampling from $\text{MGV}(\hat{\mu}_i, \hat{\varsigma}_i)$ and back-transforming to get the

106 empirical distribution of age proportions.

107 Abundance $Y_i$ taken at different locations is considered to be spatially correlated. However, spatial patterns

108 may be blurred by factors affecting abundance observations unrelated to population size such as lighting

109 and sea conditions (Petrakis et al. 2001; Chen et al. 2004; Hjellvik et al. 2004; Johnsen and Lilende 2007).

110 If such factors are also measured, a generalised linear model (GLM) (McCullagh and Nelder 1991) can be

111 used to estimate their effects and calibrate the observations by predicting to similar hauling conditions. This

112 *calibrated abundance* data is computed by using the GLM to predict yearly abundance in specific conditions,

113 the reference conditions and adding the deviance residuals. GLMs applied at this stage are also able to

114 deal with asymmetry and over-dispersion caused by the large number of null catches (Martin et al. 2005;

115 Maunder and Punt 2004) or the occurrence of very large catches (Smith 1997; Kappenman 1999).

116 Consider now the calibrated abundance $Z_i(x_k)$, in year $i$ at location $x_k$ where $k = 1, \ldots, K$ indexes sampled

locations in the study region $A \subset \mathbb{R}^2$. We model $Z_i(x_k)$ with a Gaussian spatial geostatistical process Diggle and Ribeiro (2007). The vector of variables $Z(x)$ can be written as $Z(x) = S(x) + \epsilon$ where $S(x)$ is a stationary Gaussian process at locations $x$, with $E[S(x)] = \beta$, $Var[S(x)] = \sigma^2$ and an isotropic correlation function $\rho(h) = Corr[S(x), S(x')]$, where $h = \|x - x'\|$ is the Euclidean distance between locations $x$ and $x'$. The terms $\epsilon$ are assumed to be mutually independent and identically distributed $\text{Gau}(0, \tau^2)$. Under these settings $Z(x) \sim \text{MVG}(\beta, \Theta)$ with $\Theta$ parametrised by $(\sigma^2, \phi, \tau^2)$, where $\phi$ is the parameter reflecting the extent of the spatial correlation. Several geostatistical methods are available to make inference about $\Theta$ (Isaaks and Srivastava 1989; Cressie 1993; Diggle et al. 1998; Chilès and Delfiner 1999; Rivoirard et al. 2000; Diggle and Ribeiro 2007). We adopt Bayesian methods to compute the posterior distributions of the correlation parameters and predictive distributions for the values of $Z(x_0)$, where $x_0$ is a grid of unsampled locations over the study area (Diggle and Ribeiro 2007). Our main goal with this approach is to take into account explicitly parameter uncertainty. Notice that $\beta$ reflects $Z(x)$ mean abundance over the study area and its posterior distribution is used to obtain the empirical distribution of $Y$ directly or back transforming if necessary. On the other hand, the predicted $Z(x_0)$ over the study area reflects the spatial distributions of abundance allowing the study of spatial patterns and their evolution by year.

The analysis $Y_i$ and $\mathbf{P}_i$ can be performed in parallel and the Monte Carlo simulations are combined to produce the distribution of abundance at age by $I_{ijs} = Y_{is}P_{ijs}$ where $s = 1, \ldots, S$ indexes simulations. Figure 2 shows the algorithm used clearly identifying the two submodels, the data used for each, how the distinct analysis progress to estimate parameters and run Monte Carlo simulations, and the final combination of both submodels into abundance at age. Statistics of interest are computed based on $I_{ijs}$ and the abundance at age simulations can be used as input to large simulation frameworks, like those requested by MSE.

All analysis were carried out using R (R Development Core Team 2007) and the add-on package geoR (Ribeiro Jr and Diggle 2001).

# 4    Results

We have started the analysis with diagnostics for the model assumptions and suitable transformations. A multinomial model without covariates was compared to another fit with age proportions explained by the total catch. The latter did not improve the fit supporting the assumption of independence between total abundance and age proportions. For the additive log-ratio transformation it is necessary to choose a reference age class and a constant to be added to the data in case of the occurrence of zero counts. Choices for age class two and a value 0.1 for the constant ensured better properties in terms of skewness and normality at transformed scale, all together inducing only a small average change rate for all ages, except for age 5 with rates up to 3, mainly due to the small values observed.

<sub>149</sub> Figure 3 shows the age compositions per year with quantile based intervals obtained from 1000 bootstrap

<sub>150</sub> simulations. The survey catchability shows a dome shape with maximums at ages 1 and 2 that present the

<sub>151</sub> highest relative catches. Shifts between ages can reflect shifts in abundance at age but can also be due to

<sub>152</sub> ageing errors, not uncommon for hake (de Pontual et al. 2006; Pineiro et al. 2007).

<sub>153</sub> Abundance observations showed greater variability than predicted by a Poisson model and a negative bino-

<sub>154</sub> mial GLM with log link function provided a better fit. The measured covariates were *dayperiod*, *fortnight*,

<sub>155</sub> *bottom salinity* and *bottom temperature*. Dayperiod aimed to capture the effect of daylight with tree lev-

<sub>156</sub> els, until one hour after sunrise, after one hour before sunset and between both limits. Fortnight captured

<sub>157</sub> seasonal effects with seven levels, from the second half of September to the end of December. Bottom tem-

<sub>158</sub> perature and salinity were included as continuous variables to capture geophysical effects. The GLM was

<sub>159</sub> fitted by firstly including and fixing the *year* effect and then testing for all the other covariates including

<sub>160</sub> second degree interactions. The analysis showed significant effects only for year, fortnight and their interac-

<sub>161</sub> tion. The non-significance of the other covariates can be explained by the fact that all hauls are executed

<sub>162</sub> with some daylight and the bottom temperature and salinity are roughly constant at the depths where most

<sub>163</sub> sampling took place. The adjusted model reduced the residual deviance in 13% which, although low, is not

<sub>164</sub> unusual for this kind of analysis (Maunder and Punt 2004).

<sub>165</sub> The calibrated data set $Z_i(x_k)$ used in the geostatistical analysis was obtained by predicting abundance per

<sub>166</sub> year for the second fortnight of October and adding these values to the corresponding deviance residuals.

<sub>167</sub> To verify the univariate normality of $Z_i(x_k)$ the Shapiro-Wilks normality test was computed and 16 out of

<sub>168</sub> 18 data sets did not reject the null hypothesis of normality at an $\alpha = 0.01$, whereas for the log-transformed

<sub>169</sub> original data set, the null hypothesis was not rejected only for one out of 18.

<sub>170</sub> Geostatistical analysis adopted the exponential correlation function with algebraic form $\rho(h) = \exp\{-h/\phi\}$

<sub>171</sub> with $\rho(h) \simeq 0.05$ for the *practical range* $h = 3\phi$. Taking into account the small data set available and the

<sub>172</sub> lack of observations at short distances, we avoid estimating any other correlation parameter from the data

<sub>173</sub> by trying to fit different correlation models. Before proceeding with inference and prediction we checked for

<sub>174</sub> anisotropy effects using profiled likelihoods (Diggle and Ribeiro 2007). The profiles obtained were too flat to

<sub>175</sub> identify anisotropy parameters and the analysis proceeded assuming an isotropic spatial process. In practice,

<sub>176</sub> anisotropy effects are extremely difficult to identify and usually require subjective information and/or a fairly

<sub>177</sub> large amount of samples which is uncommon on bottom trawl surveys data sets. Considering isotropy and

<sub>178</sub> the small number of samples available per year, we rotated the southern continental shelf 90º clockwise, so

<sub>179</sub> that it became aligned with the western coast, in order to use as much information as possible for inference

<sub>180</sub> on model parameters.

<sub>181</sub> The priors for the correlation parameters were set based on our experience modelling this data (Mendes 2007;

Jardim and Ribeiro Jr. 2007, in press) and our knowledge of the stochastic process correlation structure. For the range parameter $\phi$ we used an exponential prior distribution with an expected value of 20km, reflecting higher beliefs on short correlations. The nugget variance parameter $\tau^2$ was reparameterized into a relative nugget $\tau^2_{REL} = \tau^2 \sigma^{-2}$ and the prior set as a zero inflated Poisson (ZIP) distribution with mean of the positive values equals to 1.25 and a probability of zero value equals to 0.25. These probabilities were initially computed for values 0 to 8 and reassigned to 9 even intervals between 0 and 2. Our choice is based on the prior belief that the GLM analysis should have removed most of the random noise from the data and $\tau^2$ is *a priori* expected to be small. On the other hand, to estimate $\tau^2$ it is necessary to have observations at the same location or at very close distances, which is operationally not feasible for BTS. For the mean parameter $\beta$ we used a flat prior. The same priors were adopted for all years. Prior and posterior distributions of $\phi$ and $\tau^2_{REL}$ are shown in Figure 4. The posterior distributions of $\phi$ showed modes approximately between 10 and 20 km, reflecting a practical correlation range between 30 and 60 km, perfectly acceptable considering the length of the Portuguese coast. For $\tau^2_{REL}$ it is clear that the data does not contain much information about the parameter and the posterior distributions are very similar to the priors, in particular in 1990 and between 1992 and 1997. This impacts prediction variances as $\tau^2$ reflects the random variability of the process.

Yearly abundance simulations were computed by $Y_{is} = \exp(\beta_{is})$ where $\beta_{is}$ are the yearly simulations of the posterior distribution of $\beta$. The requirement to back transform $\beta_{is}$ was caused by the log transformation used to compute the calibrated abundance with the GLM. The abundance index and the 95% credibility intervals were obtained computing the median and the 0.025 and 0.975 quantiles of $Y_i$ (Figure 5). Abundances showed a cyclic pattern with high values in 1991, 1997, 2001 and 2005; and low values in 1993, 1996, 1999, 2003 and 2006. There is a persistent increase from 1993 although still within the historical limits. The credibility intervals are asymmetric and showed larger intervals in the highest estimates as expected by the GLM log transformation. Table 1 presents several metrics computed using design statistics and geostatistics. Considering the asymmetry of $Y_{is}$ we computed the relative median absolute deviation, the ratio between the median absolute deviation and the median, that can be seem as a robust adimensional indicator of precision, comparable to the coefficient of variation. The values obtained by geostatistics are lower than those obtained by design statistics. This result can be explained by a screening effect (Isaaks and Srivastava 1989) that downweights groups of observations nearby as the information contained in each observation becomes redundant. In such cases aggregations of high observations in space (Figure 1) have a lower impact on the results of the geostatistical analysis than on design-based methods given the sensibility of the sample mean to high values. The higher precision obtained with design estimators is apparently over-optimistic for BTS, where sample sizes are always small due to the operational costs. The amount of information contained in each sample is overestimated when ignoring the correlation between samples, leading to an

7

underestimated variance. Geostatistical results present a relative median absolute deviation between 14 and 25, in agreement with other studies (*e.g.* see Smith and Gavaris 1993; Dingsor 2005; Sousa et al. 2007; Roa-Ureta and Niklitschek in press).

Spatial predictions were carried out on a grid over the study area with locations at 5 km of each other resulting in 1255 locations within the study area. Figure 6 presents the spatial distribution of hake over the study area standardised by the maximum in each year so that the year effect was removed highlighting the spatial effect present on the maps. It is possible to identify persistent areas of high abundance on the west coast at latitudes approximately of 4150km (UTM), 4280km (UTM) and 4400km (UTM). The first and second areas are known recruitment spots and the last one is less persistent, but also known to be an area of high recruitment.

Abundance at age and year are presented in the bottom panel of Table 2 with the relative median absolute deviation between brackets. As with $Y_i$ the estimates of abundance at age are lower and less precise than the design-based ones, resulting from the fact that $I_{ij}$ accounts for the variability of both, $Y_i$ and $\mathbf{P}_i$. The same reasoning regarding the screening effect and variance underestimation also applies here. A comparison between design-based statistics and our estimates is presented in Figure 7, with both time series standardised to zero mean and unit variance. In general both series are similar identifying the same maxima and minima, the highest differences arise in ages 4 and 5 which are not well represented on the survey catches.

# 5 Discussion

The model proposed considers that modelling abundance at age requires two main characteristics to be taken into account, the aggregation of individuals of similar length and the spatial patterns of abundance, accounting for the major sources of variability. The separation of the age compositions from the age-aggregated abundance allows suitable models to be applied to each variable, improving the analysis and increasing the flexibility of the model. Age structures were studied by compositional data analysis considering the full covariance structure of age compositions. Age-aggregated data was modelled with geostatistical methods explicitly taking into account the correlation between abundance at different locations. Geostatistical models for compositional data (Walvoort and de Gruijter 2001; Pawlowsky-Glahn and Olea 2004) are still in development and our view is that the scarcity of data provided by BTS tend to impair the use of data demanding approaches.

An important feature of the proposed model is its full parametric specification allowing for the usage of Monte Carlo simulation methods, providing ways to overcome difficulties in obtaining an analytical expression for the statistical distribution of abundance at age, while still allowing for the computation of several statistics of interest. Outputs can also be used as inputs for larger simulation frameworks like MSE. MSE constitutes

8

a modern and sophisticated approach to management of fisheries and ecosystems but, despite its formal complexity, the output and advice obtained is equally reliant on the quality of its inputs. The approach presented in this work is one step forward on providing stochastic input parameters. Additionally the methods advocated in this paper produce several abundance indicators that provide an overview of abundance along different perspectives. The analysis of age compositions provides an insight on how the population structure evolves over time. The geostatistical submodel returns abundance indicators in both space and time perspectives, whereas the possibilities of explicitly modelling space-time interactions can be investigated (Silva et al., 2007).

In practice, modelling abundance data requires several adjustments depending on the species, area and main objectives. Our case study allowed us to point out possible solutions but it will always be necessary to find appropriate solutions considering individual characteristics of the problem at hand. The application presented assumed that age compositions were independent from age aggregated catches, an assumption supported by the exploratory data analysis. In more general terms this issue can be solved by post-stratification of the study area into strata where this assumption stands, either discretizing the age aggregated catches and modelling each data set independently or explicitly modelling this relation.

The problem of asymmetry and over-dispersion surfaced during the analysis of our data set, caused by a large number of null or small observations and occasional very large catches. The GLM with negative binomial errors used to calibrate the observations provides a way to sort out such problems, and explained a considerable part of the spatially unstructured variability, as indicated by the low values of $\tau^2$. On the other hand, the problem of modelling zero observations is restricted to $\mathbf{P}_i$ and had a negligible impact on the geostatistical analysis which uses the age-aggregated catches, less likely to have null observations. This is another advantage of the proposed approach, as modelling abundance at age using geostatistics can be severely limited by zero values, commonly present on ages poorly represented in the sample. Attempts to apply geostatistical models separately to different ages will most likely result in different and eventually conflicting inferences on the correlation parameters, and inconsistent spatial predictions.

# 6   Acknowledgements

9

# References

Aitchison, J., 1982. The statistical analysis of compositional data. Journal of the Royal Statistical Society. Series B 44 (2), 139–177.

Aitchison, J., 2003. The Statistical Analysis of Compositional Data. The Blackburn Press.

Anonymous, 2002. Report of the International Bottom Trawl Survey Working Group. CM D:03, ICES.

Babak, O., Hrafnkelsson, B., Palsson, O., 2007. Estimation of the length distribution of marine populations in the gaussian-multinomial settings using the methods of moments. Journal of Applied Statistics 34 (8), 985–995.

Brynjarsdottir, J., Stefansson, G., 2004. Analysis of cod catch data from icelandic groundfish surveys using generalized linear models. Fisheries Research 70, 195–208.

Chen, J., Thompson, M., Wu, C., 2004. Estimation of fish abundance indices based on scientific research trawl surveys. Biometrics 60, 116–123.

Chilès, J.-P., Delfiner, P., 1999. Geostatistics: Modeling Spatial Uncertainty. Wiley, New York.

Cochran, W., 1960. Sampling Techniques. Statistics. John Wiley and Sons, INC, New York.

Cressie, N., 1993. Statistics for spatial data - Revised Edition. John Wiley and Sons, New York.

de Pontual, H., Groison, A. L., Pineiro, C., Bertignac, M., 2006. Evidence of underestimation of european hake growth in the bay of biscay, and its relationship with bias in the agreed method of age estimation. ICES J. Mar. Sci. 63 (9), 1674–1681.

Diggle, P., Ribeiro, P., 2007. Model-based Geostatistics. Springer, New York.

Diggle, P. J., Tawn, J. A., Moyeed, R. A., 1998. Model-based geostatistics (with discussion). Appl. Statist. 47, 299–350.

Dingsor, G. E., 2005. Estimating abundance indices from the internationsl 0-group fish survey in the Barents Sea. Fisheries Research (72), 205–218.

Efron, B., Tibshirani, R., 1993. An Introduction to the Bootstrap. Chapman & Hall, New York.

Hammond, P. S., Donovan, G. P., in press. Development of the iwc revised management procedure. Journal of Cetacean Research and Management, Special Edition 4.

Hjellvik, V., Godo, O., Tjostheim, D., 2004. Decomposing and explaining the variability of bottom trawl survey data from the barents sea. Sarsia: North Atlantic Marine Science 89 (15), 196–210.

Hrafnkelsson, B., Stefansson, G., 2004. A model for categorical length data from groundfish surveys. Canadian Journal of Fisheries and Aquatic Science 61, 1135–1142.

Isaaks, E., Srivastava, M., 1989. An Introduction to Applied Geostatistics. Oxford University Press, New York.

Jardim, E., Ribeiro Jr., P., in press. Geostatistical tools for assessing sampling designs applied to a portuguese bottom trawl survey field experience. Scientia Marina.

Jardim, E., Ribeiro Jr., P. J., 2007. Geostatistical assessment of sampling designs for portuguese bottom trawl surveys. Fisheries Research 85, 239–247.

Johnsen, E., Lilende, T., 2007. Factors affecting the diel variation in commercial cpue of namibian hake - can new information improve standard survey estimates ? Fisheries Research 88, 70–79.

Johnston, S. J., Butterworth, D. S., 2005. The evolution of operational management procedures for the south african west coast rock lobster fishery. New Zealand Journal of Marine and Freshwater Research 39, 687–702.

Kappenman, R. F., 1999. Trawl survey based abundance estimation using data sets with unusually large catches. ICES Journal of Marine Science 56, 28–35.

Kell, L. T., Mosqueira, I., Grosjean, P., Fromentin, J.-M., Garcia, D., Hillary, R., Jardim, E., Mardle, S., Pastoors, M. A., Poos, J. J., Scott, F., Scott, R. D., 2007. Flr: an open-source framework for the evaluation and development of management strategies. ICES Journal of Marine Science 64 (4), 640–646.
URL http://icesjms.oxfordjournals.org/cgi/content/abstract/64/4/640

Martin, T. G., Wintle, B. A., Rhodes, J. R., Kuhnert, P. M., Field, S. A., Low-Choy, S. J., Tyre, A. J., Possingham, H. P., 2005. Zero tolerance ecology: improving ecological inference by modelling the source of zero observations. Ecology Letters 8 (11), 1235–1246.

Maunder, M., Punt, A., 2004. Standardizing catch and effort data: a review of recent approaches. Fisheries Research 70, 141–159.

McConnaughey, R. A., Conquest, L. L., 1993. Trawl survey estimation using a comparative approach based on lognormal theory. Fishery Bulletin 91 (1), 107–118.

McCullagh, P., Nelder, J., 1991. Generalized Linear Models, 2nd Edition. No. 37 in Monographs on Statistics and Applied Probability. Chapman and Hall.

334  Mendes, K. F. T. J. M., 2007. A bayesian hierarchical model for over-dispersed count data: a case study for

335  abundance of hake recruits. Environmetrics 18, 27–53.

336  URL http://dx.doi.org/10.1002/env.800

337  Murteira, B. J., 1990. Probabilidades e Estatistica, 2nd Edition. Vol. 2. McGraw-Hill.

338  O'Neill, M. F., Faddy, M. J., 2003. Use of binary and truncated negative binomial modelling in the analysis

339  of recreational catch data. Fisheries Research 60, 471–477.

340  Pawlowsky-Glahn, V., Olea, R. A., 2004. Geostatistical Analysis of Compositional Data. Oxford University

341  Press, USA.

342  Pennington, M., 1983. Efficient estimators of abundance, for fish and plankton surveys. Biometrics 39 (1),

343  281–286.

344  Petrakis, G., MacLennan, D., Newton, A., 2001. Day-night and depth effects on catch rates during trawl

345  surveys in the North Sea. ICES Journal of Marine Science (58), 50–60.

346  Piet, G., 2002. Using external information and gams to improve catch-at-age indices for north sea plaice and

347  sole. ICES Journal of Marine Science 59, 624–632.

348  Pineiro, C., Rey, J., de Pontual, H., Goni, R., 2007. Tag and recapture of european hake (merluccius

349  merluccius l.) off the northwest iberian peninsula: First results support fast growth hypothesis. Fisheries

350  Research 88, 150–154.

351  URL http://www.sciencedirect.com/science/article/B6T6N-4PK7P0J-1/2/ff9effeb42b316b491cbb28f83d68fd2

352  Pradhan, N. C., Leung, P., 2006. A poisson and negative binomial regression model of sea turtle interactions

353  in hawaii's longline fishery. Fisheries Research 78, 309–322.

354  Punt, A. E., Pribac, F., Taylor, B., Walker, T., 2005. Harvest strategy evaluation for school and gummy

355  shark. Journal of Northwest Atlantic Fisheries Science (Online) 35.

356  R Development Core Team, 2007. R: A Language and Environment for Statistical Computing. R Foundation

357  for Statistical Computing, Vienna, Austria, ISBN 3-900051-07-0.

358  URL http://www.R-project.org

359  Ribeiro Jr, P. J., Diggle, P. J., June 2001. geoR: a package for geostatistical analysis. R-NEWS 1 (2), 14–18,

360  iSSN 1609-3631.

361  URL http://CRAN.R-project.org/doc/Rnews/

362  Rivoirard, J., Simmonds, J., Foote, K., Fernandes, P., Bez, N., 2000. Geostatistics for Estimating Fish

363  Abundance. Blackwell Science, London, England.

364 Roa-Ureta, R., Niklitschek, E., in press. Biomass estimation from surveys with likelihood-based geostatistics.
365    ICES Journal of Marine Science, 12.

366 Silva, A. S., Ribeiro Jr, P. J., Elmatzoglou, I., 2007. Modelagem geoestatï¿œstica utilizando a famï¿œelia de
367    gneiting de funï¿œeï¿œees de covariï¿œencia espaï¿œeo-temporais. Revista de Matemï¿œetica e Estatï¿œestica,
368    v. 25, p. 65-83, 2007. 25 (2), 65–83.

369 Smith, S., Gavaris, S., 1993. Improving the Precision of Abundance Estimates of Eastern Scotian Shelf
370    Atlantic Cod from Bottom Trawl Surveys. North American Journal of Fisheries Mangement (13), 35–47.

371 Smith, S. J., 1988. Evaluating the efficiency of the delta-distribution mean estimator. Biometrics 44 (2),
372    485–493.

373 Smith, S. J., 1990. Use of statistical models for the estimation of abundance from groundfish trawl survey
374    data. Canadian Journal of Fisheries and Aquatic Science 47, 894–903.

375 Smith, S. J., 1997. Bootstrap confidence limits for groundfish trawl survey estimates of mean abundance.
376    Canadian Journal of Fisheries and Aquatic Science 54, 616–630.

377 Sousa, P., Lemos, R., Gomes, M., Azevedo, M., 2007. Analysis of horse mackerel, blue whiting, and hake catch
378    data from portuguese surveys (1989-1999) using an integrated glm approach. Aquatic Living Resources
379    20, 105–116.

380 Stefansson, G., 1996. Analysis of groundfish survey abundance data: combining the glm and delta approaches.
381    ICES Journal of Marine Science 53, 577–588.

382 Thompson, S., 1992. Sampling. Statistics. John Wiley & Sons, INC, New York.

383 Walvoort, de Gruijter, 2001. Compositional kriging: A spatial interpolation method for compositional data.
384    Mathematical Geology 33, 951–966.
385    URL http://dx.doi.org/10.1023/A:1012250107121

Table 1: Age aggregated abundance estimates by design statistics and geostatistics. The design statistics were the stratified mean, $\hat{Y}$, its standard deviation, $\sigma_{\hat{Y}}$, and coefficient of variation, $\text{CV}_{\hat{Y}}$. The geostatistics were the median, $\tilde{Y}$, the median absolute deviation, $\text{MAD}_{\tilde{Y}}$, the relative median absolute deviation, $\text{RMAD}_{\tilde{Y}}$, the 0.025, $\text{Q}(\tilde{Y}, 0.025)$, the 0.975 percentiles, $\text{Q}(\tilde{Y}, 0.975)$, and the interquartile range, $\text{IQR}_{\tilde{Y}}$.

| | | design statistics | | | geostatistics | | | | | |
|------|-------|-------------|--------------------|-------------------|-------------|----------------------|-----------------------|---------------------------|---------------------------|-----------------------|
| Year | hauls | $\hat{Y}$ | $\sigma_{\hat{Y}}$ | $\text{CV}_{\hat{Y}}$ | $\tilde{Y}$ | $\text{MAD}_{\tilde{Y}}$ | $\text{RMAD}_{\tilde{Y}}$ | $\text{Q}(\tilde{Y}, 0.025)$ | $\text{Q}(\tilde{Y}, 0.975)$ | $\text{IQR}_{\tilde{Y}}$ |
| 1989 | 130 | 59.2 | 1.7 | 0.03 | 33.6 | 6.6 | 0.2 | 21.2 | 49.7 | 28.4 |
| 1990 | 108 | 157 | 9.7 | 0.06 | 38.9 | 6.4 | 0.16 | 25.9 | 52.8 | 26.9 |
| 1991 | 80 | 194.1 | 12.2 | 0.06 | 154.8 | 27.4 | 0.18 | 101.3 | 250.4 | 149.1 |
| 1992 | 44 | 65.3 | 3.2 | 0.05 | 46.1 | 10.4 | 0.22 | 26.4 | 79.5 | 53 |
| 1993 | 58 | 54.1 | 4.5 | 0.08 | 8.1 | 1.5 | 0.18 | 5.5 | 11.9 | 6.5 |
| 1994 | 76 | 95.9 | 4.7 | 0.05 | 61.8 | 8.5 | 0.14 | 46.6 | 82.3 | 35.7 |
| 1995 | 80 | 85.2 | 4.1 | 0.05 | 59.4 | 8.5 | 0.14 | 42.1 | 80.7 | 38.5 |
| 1996 | 63 | 44.6 | 2.3 | 0.05 | 25.1 | 6.4 | 0.25 | 15.7 | 44.1 | 28.4 |
| 1997 | 51 | 207.2 | 21.5 | 0.1 | 123.9 | 20.1 | 0.16 | 86.9 | 188.4 | 101.4 |
| 1998 | 64 | 139.8 | 7.8 | 0.06 | 109.4 | 21.3 | 0.19 | 65.5 | 164.5 | 99 |
| 1999 | 71 | 71.2 | 2.5 | 0.04 | 27.3 | 5.8 | 0.21 | 16.1 | 42.2 | 26.1 |
| 2000 | 65 | 102.2 | 5.8 | 0.06 | 89.2 | 14.3 | 0.16 | 63 | 134.3 | 71.4 |
| 2001 | 58 | 164 | 15.3 | 0.09 | 140.3 | 23.2 | 0.17 | 91 | 199 | 107.9 |
| 2002 | 66 | 117.5 | 7.9 | 0.07 | 75 | 18.7 | 0.25 | 41.8 | 120.4 | 78.6 |
| 2003 | 72 | 55.3 | 2 | 0.04 | 41.5 | 8.4 | 0.2 | 25.6 | 65.2 | 39.6 |
| 2004 | 79 | 124.4 | 6.3 | 0.05 | 77.8 | 19.4 | 0.25 | 42.6 | 134.7 | 92.1 |
| 2005 | 87 | 214 | 9.4 | 0.04 | 153 | 29.7 | 0.19 | 93.6 | 235.2 | 141.7 |
| 2006 | 88 | 125.9 | 4.4 | 0.03 | 42.6 | 8.8 | 0.21 | 26.4 | 66.3 | 39.9 |

Table 2: Abundance at age estimates by design statistics on the top panel and this study on the bottom panel. The design statistics are the stratified mean and between brackets its coefficient of variation. The estimates provided by this study are the median and between brackets the relative median absolute deviation.

| Estimator | Year | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|---|
| Design based | 1989 | 12.9 (0.08) | 20.1 (0.05) | 16.9 (0.04) | 7.4 (0.06) | 1.5 (0.09) | 0.4 (0.14) |
| | 1990 | 82.1 (0.11) | 45.4 (0.05) | 19.3 (0.05) | 7.4 (0.05) | 2.4 (0.07) | 0.4 (0.12) |
| | 1991 | 56.6 (0.14) | 82.4 (0.10) | 36.7 (0.11) | 14.6 (0.08) | 3.1 (0.09) | 0.6 (0.12) |
| | 1992 | 12.1 (0.16) | 20.4 (0.09) | 19.3 (0.08) | 10.2 (0.07) | 2.7 (0.10) | 0.6 (0.17) |
| | 1993 | 23.2 (0.18) | 17.1 (0.09) | 8.6 (0.11) | 3.6 (0.10) | 1.3 (0.14) | 0.3 (0.32) |
| | 1994 | 18.5 (0.14) | 51.4 (0.07) | 18.2 (0.08) | 5.9 (0.10) | 1.5 (0.15) | 0.3 (0.21) |
| | 1995 | 2.1 (0.16) | 34.6 (0.09) | 37.2 (0.07) | 8.1 (0.13) | 2.9 (0.17) | 0.4 (0.23) |
| | 1996 | 9.0 (0.10) | 15.1 (0.09) | 10.8 (0.12) | 6.9 (0.12) | 1.9 (0.16) | 0.9 (0.17) |
| | 1997 | 40.4 (0.22) | 70.4 (0.18) | 83.7 (0.18) | 8.7 (0.17) | 2.3 (0.29) | 1.6 (0.32) |
| | 1998 | 54.0 (0.11) | 46.5 (0.10) | 22.8 (0.08) | 12.3 (0.09) | 3.0 (0.13) | 1.1 (0.17) |
| | 1999 | 9.1 (0.12) | 26.9 (0.05) | 25.0 (0.07) | 7.8 (0.09) | 2.0 (0.13) | 0.4 (0.22) |
| | 2000 | 29.9 (0.14) | 39.3 (0.09) | 21.4 (0.08) | 8.9 (0.10) | 1.7 (0.12) | 1.0 (0.16) |
| | 2001 | 50.9 (0.23) | 73.9 (0.13) | 22.2 (0.10) | 14.3 (0.09) | 2.1 (0.15) | 0.6 (0.20) |
| | 2002 | 43.5 (0.16) | 37.1 (0.09) | 26.8 (0.08) | 7.5 (0.11) | 2.1 (0.15) | 0.4 (0.26) |
| | 2003 | 5.9 (0.08) | 28.6 (0.05) | 13.2 (0.08) | 6.1 (0.09) | 1.3 (0.15) | 0.2 (0.27) |
| | 2004 | 42.5 (0.10) | 48.6 (0.08) | 22.8 (0.08) | 7.9 (0.11) | 1.7 (0.16) | 0.8 (0.18) |
| | 2005 | 105.8 (0.08) | 67.5 (0.05) | 30.2 (0.06) | 7.8 (0.10) | 2.0 (0.13) | 0.7 (0.20) |
| | 2006 | 44.7 (0.07) | 35.4 (0.06) | 32.6 (0.06) | 10.0 (0.09) | 2.5 (0.13) | 0.6 (0.21) |
| This study | 1989 | 2.9 (0.25) | 9.8 (0.21) | 12.2 (0.20) | 6.4 (0.22) | 1.6 (0.24) | 0.7 (0.25) |
| | 1990 | 3.9 (0.26) | 13.6 (0.20) | 11.9 (0.19) | 6.0 (0.23) | 2.4 (0.24) | 0.7 (0.25) |
| | 1991 | 14.8 (0.32) | 51.3 (0.25) | 52.0 (0.23) | 25.5 (0.26) | 7.0 (0.30) | 2.0 (0.30) |
| | 1992 | 2.7 (0.40) | 9.1 (0.31) | 13.5 (0.27) | 13.8 (0.26) | 4.7 (0.34) | 1.5 (0.38) |
| | 1993 | 1.2 (0.30) | 2.6 (0.24) | 2.2 (0.23) | 1.2 (0.29) | 0.5 (0.29) | 0.2 (0.33) |
| | 1994 | 5.2 (0.24) | 26.3 (0.21) | 15.3 (0.20) | 10.5 (0.23) | 3.3 (0.26) | 0.9 (0.27) |
| | 1995 | 1.0 (0.30) | 19.0 (0.19) | 27.5 (0.16) | 8.2 (0.19) | 2.8 (0.23) | 0.6 (0.26) |
| | 1996 | 2.6 (0.34) | 8.7 (0.30) | 6.4 (0.28) | 4.6 (0.28) | 1.7 (0.33) | 1.1 (0.32) |
| | 1997 | 2.9 (0.38) | 25.9 (0.29) | 78.4 (0.18) | 11.7 (0.25) | 2.5 (0.29) | 1.8 (0.31) |
| | 1998 | 16.2 (0.36) | 29.0 (0.26) | 27.5 (0.23) | 24.5 (0.26) | 6.8 (0.31) | 2.7 (0.31) |
| | 1999 | 1.7 (0.31) | 8.4 (0.26) | 12.3 (0.21) | 3.7 (0.26) | 0.7 (0.28) | 0.2 (0.30) |
| | 2000 | 7.8 (0.32) | 25.6 (0.23) | 32.8 (0.19) | 16.6 (0.22) | 3.7 (0.24) | 2.5 (0.25) |
| | 2001 | 11.7 (0.31) | 49.1 (0.25) | 42.7 (0.22) | 29.5 (0.24) | 3.8 (0.28) | 1.8 (0.29) |
| | 2002 | 12.1 (0.32) | 23.7 (0.3) | 26.8 (0.27) | 7.8 (0.29) | 2.5 (0.32) | 0.9 (0.35) |
| | 2003 | 3.6 (0.27) | 17.9 (0.24) | 12.7 (0.22) | 5.1 (0.26) | 1.4 (0.29) | 0.5 (0.28) |
| | 2004 | 15.7 (0.29) | 37.5 (0.25) | 17.1 (0.3) | 4.5 (0.33) | 1.5 (0.32) | 1.0 (0.33) |
| | 2005 | 37.2 (0.26) | 68.0 (0.21) | 33.8 (0.24) | 9.5 (0.26) | 2.5 (0.28) | 1.3 (0.29) |
| | 2006 | 5.3 (0.29) | 13.0 (0.23) | 15.9 (0.23) | 6.3 (0.24) | 1.5 (0.27) | 0.5 (0.28) |

Figure 1: Yearly maps with locations of hauls (+) and observed catches of Hake (*Merluccius merluccius*) during the Autumn series of the Portuguese bottom trawl survey. The gray circles are proportional to the logarithm of the numbers of individuals caught per hour. The full line represents the Portuguese continental coast.
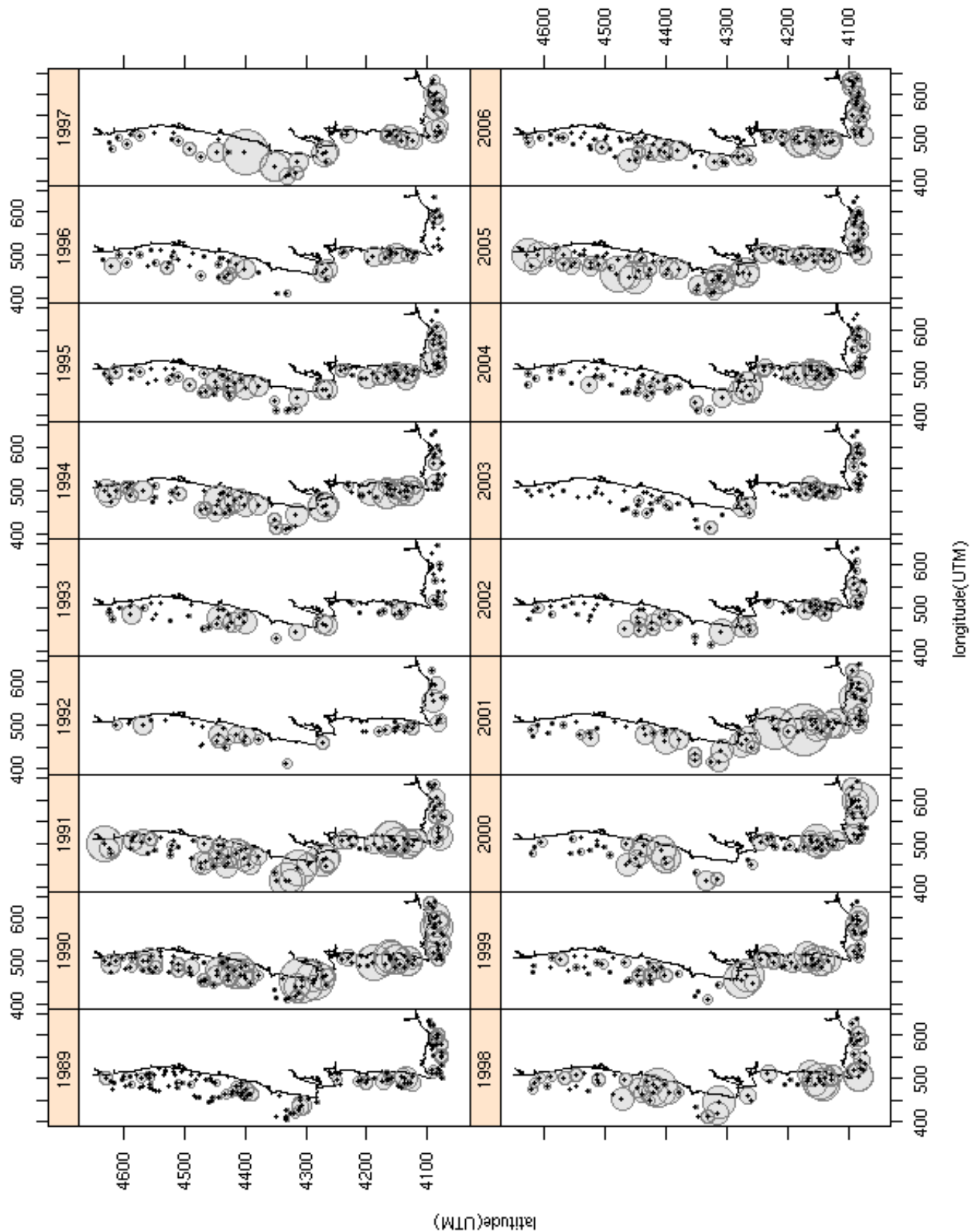
Figure 2: Graphical representation of the algorithm used for analysis showing a clear separation of yearly abundance at age, $I$, in two branches. On the the left the age structure, $P$, is analysed with compositional data analysis, and on the write the spatial distribution $Y$ is analysed with geostatistical methods. The last procedure is to combine the simulations of both variables to compute the stochastic distribution of the abundance at age per year. The round boxes represent data and the sharp boxes represent methods. $D$ is the transformed compositional data; MVG=multivariate Gaussian distribution; $\Lambda$, $\Sigma$, $\mu$ and $\varsigma$ are parameters of $D$; $Z(x)$ is a stationary spatial process; $\beta$ and $\Theta$ are parameters of the spatial models with $\sigma^2 = $ sill, $\phi = $ correlation range and $\tau^2 = $ nugget; $x_0$ is a grid of unsampled locations; $i$ indexes years, $j$ indexes ages and $s$ indexes simulations.
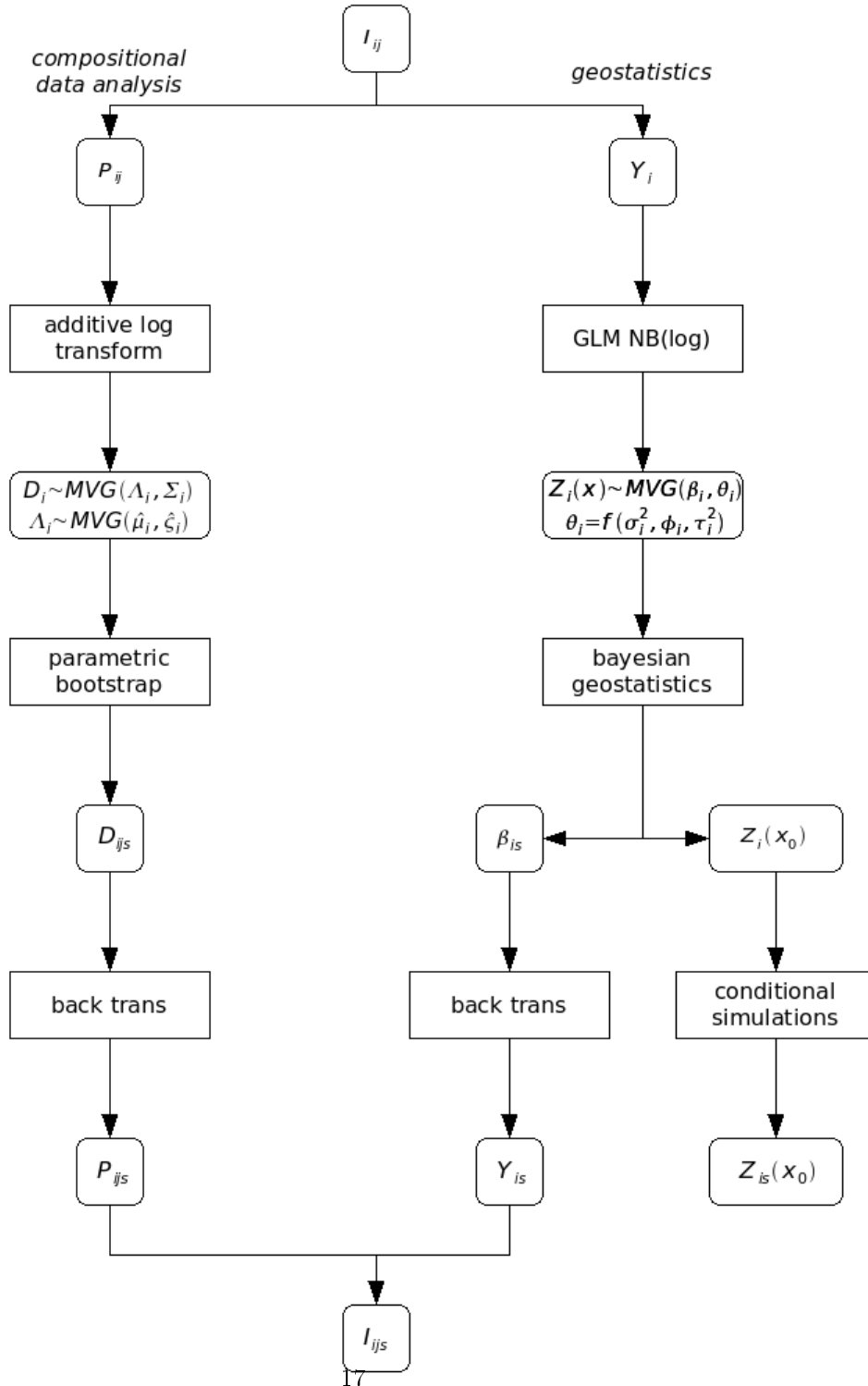
Figure 3: Age compositions empirical distribution obtained by parametric bootstrap. The full circle represents the median proportion and the gray lines represent the confidence interval computed by the 0.025 and 0.975 percentiles.
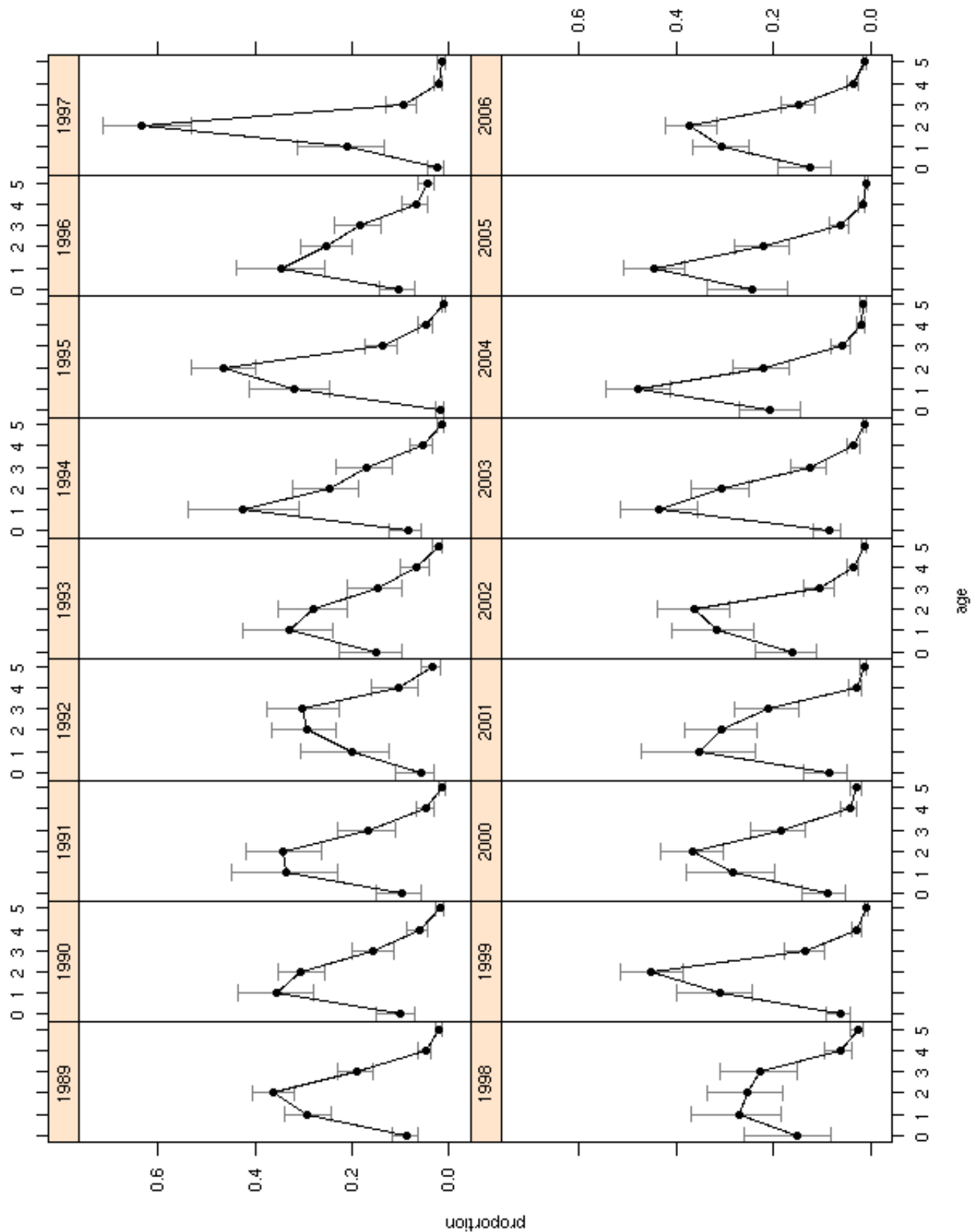
Figure 4: Yearly priors and posteriors for the correlation range $\phi$ and the relative nugget $\tau_{REL}^2$ used for the geostatistical analysis of the calibrated data set. The dashed line represents the priors for each parameter, kept constant for all data sets. The full line represents the posteriors obtained per year for each data set.
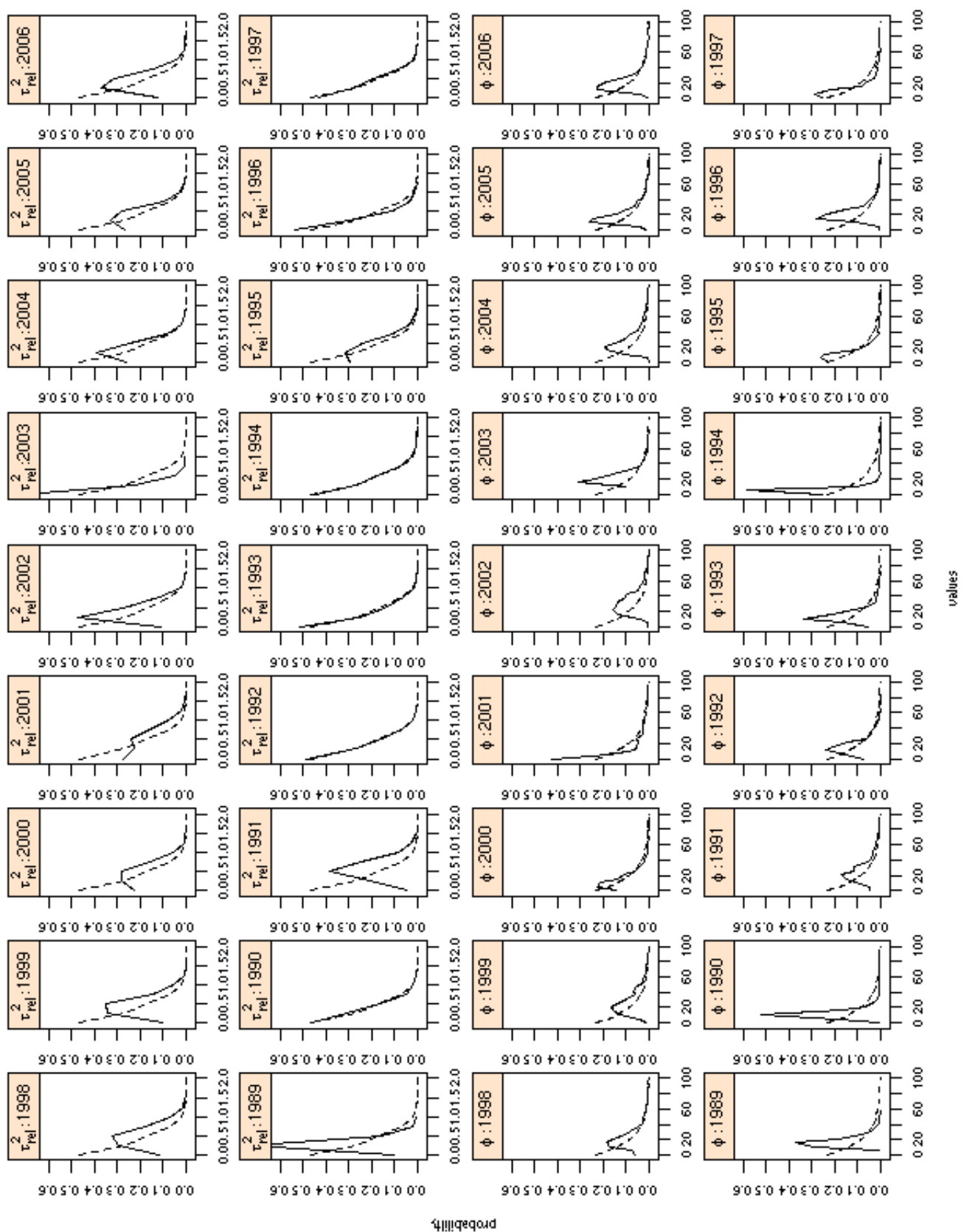
Figure 5: Yearly abundance estimates. The black circle represents the median abundance and the gray lines represent the confidence interval computed by the 0.025 and 0.975 percentiles.
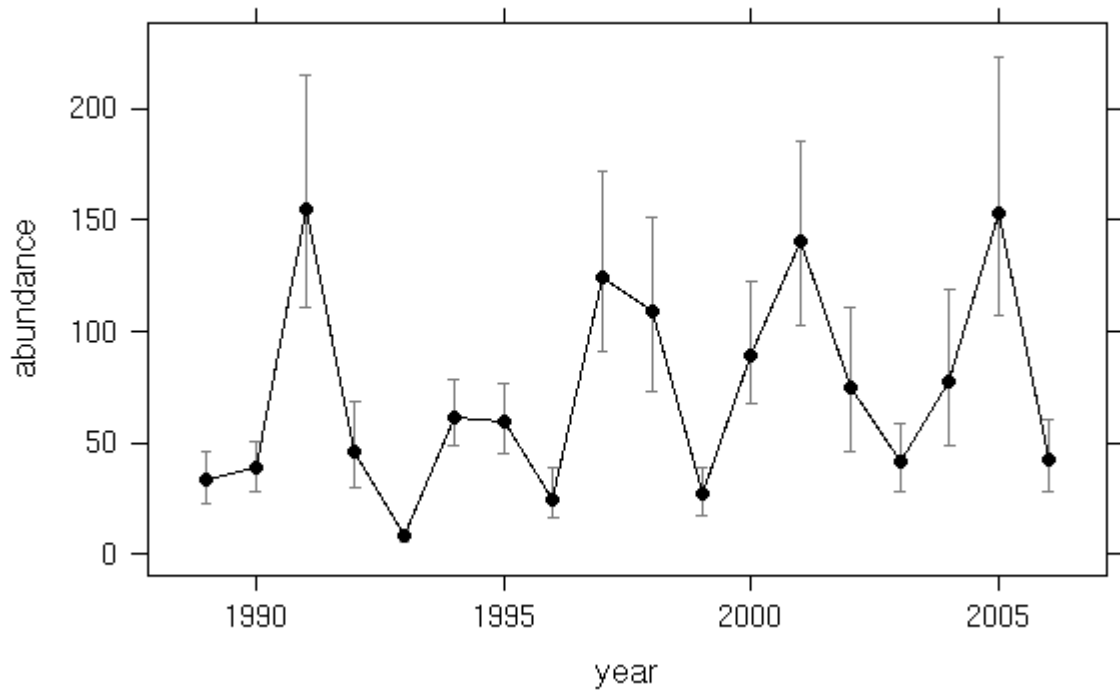
Figure 6: Spatial distribution of age aggregated abundance by year, standardised to the second fortnight of October. The gray degrees are proportional to the number of individuals caught by unit effort, rescaled to the maximum estimate within each year. The black color represent 1 and the white colour represents 0.
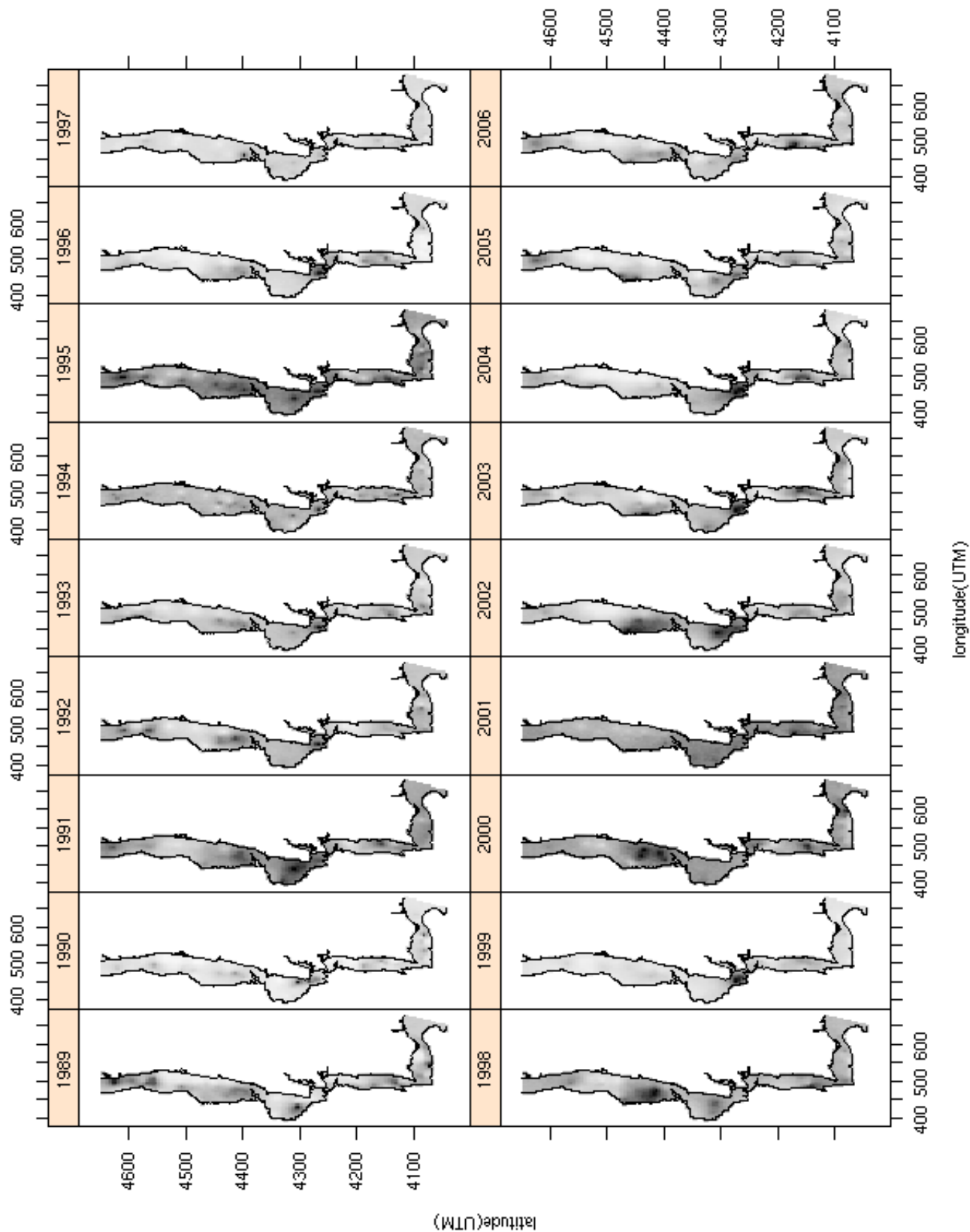
Figure 7: Abundance at age and year standardised to have mean 0 and variance 1. Design estimates in dashed line and geostatistical estimates in full line.