Not to be cited without prior reference to the authors.

# Modelling Herring Distributions in Space and Time

*Liz Clarke'. Daniel Stahl' and John* **Simmonds**[2]
*'Mathematical Institute, University of St Andrews, UK*
[2]**FRS** *Marine Laboratory Aberdeen, UK*

International Bottom Trawl survey (IBTS) data for herring between 1983-1997 are analysed using generalized additive models (GAMs; Hastie and Tibshirani, 1990) to assess the changes in spatial distribution of North Sea herring over age and time. The descriptive powers of covariates such as latitude, longitude, depth, bottom temperature, bottom salinity and median sediment grain size are explored. Model selection and the effect of smoothing on the abundance estimates is considered. A method of correcting the catches for light level and vessel effects is discussed. The resulting abundance estimates are compared to the current estimation method and ICA predictions for North Sea herring.

**Introduction**

International bottom trawl surveys (IBTS) are research surveys that have been conducted by several countries in the first quarter of the year (January-March) for many years. An even coverage is attempted, with two vessels surveying each ICES rectangle where possible. Figure 1 shows an example of the coverage of each vessel in one survey. Since 1983, efforts have been made to standardise survey protocol so that the trawls for different vessels are comparable and each trawl is a measurement of catch per unit effort (CPUE). Thus IBTS data can be used for calculating abundance indices over the North Sea. Further details of the IBTS survey can be found in Heessen *et al.* (1997). The methodology currently used by ICES to calculate the abundance indices averages catch for each ICES rectangle then sums over the rectangles in the area of interest. This corrects for oversampling within rectangles.

Exploratory data analysis indicated that the size of the catch depends on the time of day at which the trawl took place, with trawls taken at night generally having lower catches than those during the day. This effect appeared to vary smoothly throughout the day, rather than as a discontinuous jump representing night and day. The biological explanation for this is that the fish move up and down the water column according to the amount of light available. If the catches are not corrected for this effect, the night hauls performed in a survey could negatively bias the resulting abundance estimate. This in turn will affect the time series of abundance indices if the proportion of night hauls changes through time. In fact, night hauls are more prevalent in the north than the south and in the 1990s than the 1980s, so this is a potential problem. Although the IBTS protocol is standardised as much as possible over the vessels, it is thought that the different vessels may still catch at significantly different rates, once the variation due to other factors has been accounted for. If this is the case. then the. catches should be corrected for the vessel effect. One way of correcting for either of these effects is to apply a correction factor to the observed trawls before analysing the data. However, generalized additive models allow the correction factor to be calculated during the modelling process, leaving the original data unaltered, but allowing the predicted catches to be calculated as if they had been caught by one vessel or at a particular light-level. This approach requires considerable overlap between the coverage of the correcting covariate and the other covariates, e.g. if one vessel operates in an area where other vessels do not survey, we cannot separate the vessel effect from the spatial effects.

Generalized additive models are very useful for describing relationships between variables when the parametric form of the relationship is not known. They essentially consist of a sum of smooth functions of different covariates, the shape of each of these functions being estimated from the data. Inference using such model-based methods is more robust to non-representative sampling than sample-based methods, such as those used by ICES. which require a random survey design. GAMs have three major uses for fisheries data, all of which are interlinked. The first is that they can be used to describe the relationships between abundance and spatial and environmental factors. Model selection can be used to determine whether these relationships are significant. The second is that they can interpolate over areas that have not been surveyed to provide a map of the distribution of the species of interest, which may be useful in its management. Finally, abundance indices can be created for the area of interest by integrating under the predicted distribution surface. Generalized additive models were first used in to estimate abundance in a marine context by Swartzman *et al.* (1992). Swartzman *et al.* (1994,1995 & 1997) use GAMs to explain the distribution of fish in relation to environmental factors. Borchers *et al* (1997) and (Augustin *et al.* 1998) used GAM estimation methods in place of stratified sampling to estimate fish egg abundance from pelagic samples and obtained dramatic improvements in precision.

**Methods**

The data for each age-class were analysed separately. It was assumed that the distribution of the fish did not change significantly over a survey, i.e. temporal changes within a quarter were ignored. The data are very skew (Figure 2), containing a large proportion of zeroes and some extremely large values, with several over one hundred thousand. The GAMs were fitted using the statistical package S-PLUS (Mathsoft Inc.). The error distributions available in S-PLUS were not suitable to model these data directly and thus the modelling was performed in two stages. We first **modelled** the probability of catching any herring of a particular age-class using a GAM with a Binomial error distribution and the logit link. Here the response variable was 0 if no herring of that age-class were caught and 1 if some herring of that age-class were caught (i.e. a positive catch). We then reduced the data to positive catches only. The positive catch was log-transformed and **modelled** with a Gaussian error distribution and the identity link (a simple additive model).

In each stage of the modelling, the response variable was allowed to depend on a smooth function of each of the following covariates: latitude, longitude, bottom depth (btm.dpt), median sediment grain size (sediment), bottom temperature (btm.temp), bottom salinity (btm.sal), light level and vessel. The first six covariates are spatial and environmental covariates to describe the distribution of herring across the North Sea, whereas the latter two are essentially correction factors for the response. Light level had not been measured directly on the surveys and so three proxies were tested. These were time of day, sun elevation (sunel) and a day/night switch to indicate whether the trawl occurred during the day (sunrise to sunset) or at night (sunset to sunrise). Previous analysis suggested that cos(time), where time is transformed to radians in such a way that cos(time) is largest at midday, might be suitable because the catches increase as towards the middle of the day and then decrease again. Sun elevation is very similar to cos(time) but also takes position into account. Each smooth was of a single covariate. The smooths used here are splines, and are represented here by $s(x, \lambda)$, where x represents the covariate of interest and $\lambda$ represents the degrees of freedom used in the spline. Large values of $\lambda$ indicate a wiggly function whereas small values represent a very smooth function. $\lambda = 1$ represents a straight line. Hastie and Tibshirani ( 1990) suggest $\lambda = 4$ as a good starting point.

The data for each age for each yearly survey were **modelled** separately at first. This enables the relationships between covariates to change for each year. The data for each age were then combined over all the years and a smooth function of year was included in the model. This enables the relationships between covariates to remain the same over years, and is particularly useful for estimating vessel effects, the effect of the vessel remains constant over time, as we hope it would in reality.

*Model Selection*

Model selection was performed using Akaike's information criterion (AIC) (cf. Burnham & Anderson. 1998) to choose between covariates and to choose the degree of smoothing required for each covariate. AIC is a measure of the model deviance corrected for the number of parameters in the model and the model with the lowest AIC among competing models is chosen. Thus if two models with the same deviance but different numbers of parameters are compared, the model with fewer parameters will be chosen. When performing model selection, we need to treat the correction factors differently from the other covariates. This is because if we need to correct the response for light-level and vessel, we should do so before considering the spatial covariates but yet, in order to decide whether we need to correct the response, we need to take into account its spatial variation. An iterative approach was therefore taken. First a model including the six spatial and environmental covariates with smooths of four degrees of freedom was fitted to the data. Then the proxy covariate for light level was chosen by including each of the following possibilities in turn: cos(time), sunel, s(sunel, 2), s(sunel, 4) and day/night as a factor (categorical variable). This is intended to test the effect of light level having accounted for the spatial variability within the data. The model with the lowest AIC value was chosen. Having chosen the light level proxy, the effect of vessel was considered. Vessel was included as a factor within the model, and AIC was used to determine whether it improved the model fit compared to the model without vessel. Having chosen our correction factors for the response, we then performed model selection on all the covariates in the model, allowing the smoothness of the latitude and longitude to vary between 1, 2, 4 and 8 degrees of freedom, and the smoothness of the other covariates to vary between 1, 2, and 4 degrees of freedom. while the correction factors could be included or excluded from the model but their degrees of freedom could not vary. Latitude and longitude were allowed more degrees of freedom than the other covariates because they are essentially proxies for other unmeasured covariates, whereas the remaining covariates were allowed fewer degrees of freedom to obtain smoother functions that could be more easily interpreted biologically.

The above method is suitable for the data for separate years, in which there are approximately 400 trawls per year. However, combining all years leads to a data set of approximately 7000 trawls. which is too large for model selection using AIC to work appropriately – the large number of observations means that models with many degrees of freedom are selected. Therefore. the smooths of the covariates for the combined models were

chosen subjectively, based on the smoothness of the covariates selected for individual years. For example, longitude was often selected with 8 degrees of freedom in the individual models for age 1 herring and so a smooth of longitude with 8 degrees of freedom was included in the model for the combined data.

*Model Checking*
Models were checked by the following graphical methods.

A plot of deviance residuals against fitted values should be randomly scattered. Trends in residuals indicate that a trend in the observed data has not been modelled adequately, either because an unsuitable error distribution and/or link function have been used or suitable covariates have not been chosen. Plots that show an increase or decrease in absolute values of the residuals as fitted values increase indicated that an unsuitable error distribution and/or link function have been used.

A plot of the observed response against the fitted values should be closely scattered around a straight line through the origin. Consistently large variations around the line indicated that the model does not explain much variability in the data. Consistent trends away from the line indicate that some trend in the observed data has not been adequately modelled.

Spatial plots of residuals. Residuals of all sizes will be evenly scattered over the range of the data. Residuals of a particular size clustered together indicate that the model has not fitted the data very well in that area.

Plots of the smooth of each covariate, with rough confidence intervals, against the covariate itself should show narrow confidence intervals around the fitted curve. Wide confidence intervals around a wiggly curve indicate that fewer degrees of freedom could be used in the smooth.

*Abundance Estimation*
Once the models had been chosen, abundance was calculated in the following manner. A grid of points was created over the range of interest (roundfish areas 1-9 for age 1 herring and roundfish areas 1-8 for older age groups). We used a grid with nodes at a spacing of 8 minutes longitude and 4 minutes latitude, i.e. approximately every 4 nautical miles. Values of each of the covariates were assigned to each of the grid points. Surfaces of the predicted probability of a positive catch and the mean value of a positive catch could then be obtained over the grid for each year. If we assume that the size of a positive catch was independent of the chance of catching something in that trawl, these expected values can be multiplied together for each grid point, to obtain a surface for the expected catch over the area. Whether this assumption is valid is open to debate. The abundance index is then calculated by averaging the expected values over the area of interest. This makes the assumption that the area of the grid squares is constant over the area of interest. This is not the case, and whilst this area could be easily included in the calculation, we have not done so here in order to make the estimates more comparable with the current method. Time series of abundance estimates for different age-classes can then be created. To calculate an abundance estimate corrected for vessel effects, abundance was predicted over the whole area for each vessel in turn, and the mean of these estimates was then calculated. To calculate an abundance estimate corrected for light level effects, abundance was predicted for several values of Sun elevation between 6 am and 11 pm, the range over which most of the trawls are taken. The mean of these estimates was then calculated. Although the model selection described above included all covariates, a reliable spatial description of sediment grain size, bottom temperature and bottom salinity for each year were not available and so model selection was performed again excluding these three covariates, and these reduced models were used for the abundance estimation.

*Comparison with Other Estimates*
In order to assess the performance of these estimates, we would like to compare them to the true abundance of herring in the North Sea. Of course this is not known, and so we use instead estimates of abundance as carried out by the ICES Herring Assessment Working Group in 2000 (ICES 2000), which uses the Integrated Catch at Age method (ICA, Patterson & Melvin, 1996), excluding the ICES IBTS indices (referred to as the "ICA predictions"). We also compare the GAM-based estimates to the current ICES IBTS indices (referred to as the "current method").

*Cohort Strength*
These estimation methods provide abundance indices rather than absolute abundances. Thus, for comparison with other methods, we need to scale the indices to the same level. We therefore scale each series by its mean, so that each value is essentially an estimate of cohort strength. Variability in the estimates of cohort strength for the same cohort can be used as an indicator of the reliability of the *series*. For example, if age 1's are relatively plentiful in 1983, then we would expect age 2's to be relatively plentiful in 1984 and age 3's in 1985. and for them to have similar cohort strength as defined above. We therefore organise the cohort strength series into cohorts, take the standard deviation for each cohort, then take the mean of these standard deviations for all cohorts as an estimate of the reliability of the estimation method.

Results

Figure 1 shows the locations of the trawls in 1993 for each vessel. The coverage of the survey area is good overall but the coverage for each vessel is concentrated in smaller areas. However, there is overlap between at least two vessels in most cases and so the incorporation a vessel factor into the model should be possible.

*Choice of* **Model** *and Explanatory Covariates*
Figure 2 shows histograms of the catch and log-transformed catch data for 1993 age l herring. The distribution of the log-transformed data is approximately Gaussian. Figure 3 shows residual plots for age 1 herring in 1993 using three combinations of error distribution and link. It is clear that, of these, only the Gaussian errors and an identity link is suitable for these data. The spatial residuals are acceptable (Figure 4) and were not improved by the inclusion of a latitude longitude interaction in the form of a 2-dimensional "loess" function. A plot of the observed catches against fitted catches is difficult to interpret (Figure 5) but is fairly representative of the kind of plot expected from fisheries data.

A smooth of sun elevation with 4 degrees of freedom was found to be the best proxy for light level for the presence/absence model in most years. In all years a model including a light level proxy modelled the data better than a model which did not correct for light level. Choice of proxy was less obvious for the model for the positive data, and in several cases a model which did not include any light level proxy was selected. However, including sun elevation as a linear term performed best over all models and so sun elevation was selected as the proxy in these models.

Table l shows the covariates selected for the best model for separate years using stepwise selection with AIC. The spatial covariates latitude, longitude and depth are nearly always included in the model, with fairly high degrees of freedom. Bottom temperature is also useful, and was usually selected with two degrees of freedom, but sediment grain size bottom salinity were less useful and were rarely selected. Vessel was selected in the final model 67% of the time for the presence model but only 4 1% of the time for the positive model, whereas with sun elevation it was the other way around. This is interesting because in the initial stages of the model selection, when vessel was not included, sun elevation was very often selected for the presence model. This indicates some confounding between sun elevation and other variables, possibly vessel. Figure 6 shows the kinds of functions fitted to these data. The confidence bands shown on this plot reflect the uncertainty due to lack of data for some values of the covariates.

Figure 7 shows the vessel effects for age1 herring for the presence/absence and positive catches over all years. "A RG", "SC02" and "WAH3" stand out as having higher catches than the others. For ARG, this is due to the fact that ARG is the only vessel that fishes in the Skagerrak (see Figure 1), where there are high catches, in other words, there is confounding between vessel and the other covariates, For the other two vessels, this can be interpreted as a vessel effect.

*Spatial Distribution*
Generic spatial distributions for each age group are shown in Figure 8. These were obtained from the models using the combined data for all years. Age l herring tend to congregate in the shallower waters off the coasts of the Netherlands, Denmark and Germany, but by the time they are age 2, they have moved over towards Northern England and Scotland, although there are some in the middle of the North Sea. Ages 3, 4 and 5 (a plus group) congregate together, mainly in the northern part of the North Sea, although there is a small separate group in the Channel between England and France.

*Abundance Estimation and Comparison with Other Estimates*
The cohort strengths of the abundance estimates for separate years are compared with the current method and the ICA predictions in Figure 9. The GAM-based method and current methods give similar results, which are not as smooth as the ICA predictions, which is to be expected. The cohort strengths from the combined year model are smoother and are in closer agreement with the ICA predictions than those for the separate years (Figure 10). Time series of abundance for the combined data with and without sun elevation corrections show that sun elevation has most effect on small herring,which are in shallower water, and that the correction actually reduces the abundance, the opposite effect to what might be expected (Figure l la). This is because the many hauls taken at high sun elevations have now been down-weighted compared to the night hauls. However, a comparison of cohort strength shows that the effect of sun elevation has not changed the time series substantially (Figure l l b). This could largely be because year is included in the model. Similarly, although correcting for vessel increases the overall abundance estimated, the time series of cohort strength scarcely changes (Figures l2a & l2b). A comparison of the overall variability of the cohort strengths shows that, not surprisingly, the ICA predictions give the most reliable estimate of **cohort** strength (Table 2). The combined model is more reliable than the models for separate years. This is because unusual years are smoothed through, as they are in the ICA prediction. More surprising is the fact that the GAM-based models for separate years give less reliable estimates than the

current method. However, it should be borne in mind, that if we think of the current method as a model in which a parameter (the mean) is estimated for each of the 150 ICES rectangles, the number of parameters in the GAM-based model is much smaller than that in the current model. The inclusion of vessel or light level corrections do not improve the reliability of the estimates.

**Conclusions**

Generalized additive models provide a means of describing the distribution of fish using smooth functions of several covariates, with few parameters. The exact covariates and the degree of smoothing used in the model do not appear to affect the abundance estimate substantially. The catches can be corrected for differences in sampling regimes (i.e. the effects of fishing at different times of day or using different vessels) provided there is suitable coverage of data. However, for these data, although these effects are significant factors in the model, their overall affect on abundance indices is small. Bootstrap confidence intervals for the abundance estimates can be estimated for this method and work is ongoing in that direction. Finally, the model can be extended further to include all age groups with covariates to represent different ages and cohorts. This kind of model is becoming closer to an age-structured population dynamics model such as that used in ICA.

**Table 1** Model selection for the models for separate years using AIC. The numbers in the columns represent the number of times each covariate was selected in the best model for each of the 15 years. The last row gives the overall percentage of times that covariates was selected over all ages.

| Presence | lat | ion | btm.dpt | sediment | btm.temp | btm.sal | sunel | vessel |
|---|---|---|---|---|---|---|---|---|
| age 1 | 15 | 13 | 12 | 4 | 6 | 10 | 7 | 8 |
| age 2 | 15 | 12 | 11 | 8 | 13 | 10 | 9 | |
| age 3 | 15 | 13 | 14 | 6 | | 7 | 6 | |
| age 4 | 14 | 13 | 14 | 6 | 12 | 8 | 6 | |
| age 5 | 15 | 15 | 13 | 6 | 9 | 7 | 7 | 9 |
| % over all | 97% | 87% | 83% | 38% | 67% | 55% | 46% | 67% |

| Positive | lat | Ion | btm.dpt | sediment | btm.temp | btm.sal | sunel | vessel |
|---|---|---|---|---|---|---|---|---|
| age 1 | 15 | 15 | 14 | 3 | 10 | 6 | 9 | 10 |
| age 2 | 15 | 13 | 11 | 8 · | 11 | 8 | | 9 |
| age 3 | 14 | 14 | 11 | 4 | 7 | 7 | 9 | 8 |
| age 4 | 14 | 13 | 9 | 1 | 10 | 5 | 10 | 4 |
| age 5 | 13 | 10 | 11 | 5 | 10 | 3 | 9 | |
| % over all | 95% | 87% | 75% | 27% | 64% | 20% | 63% | 41% |

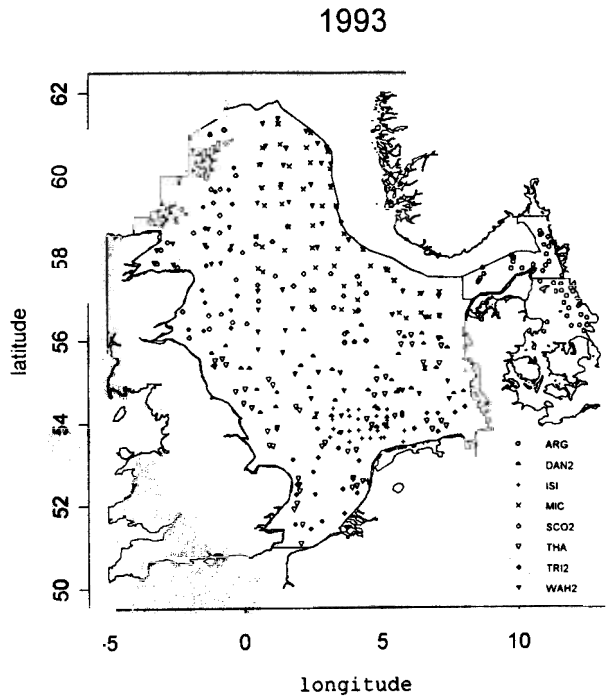**Table 2** Mean standard deviation of cohort strength for the different estimation methods.

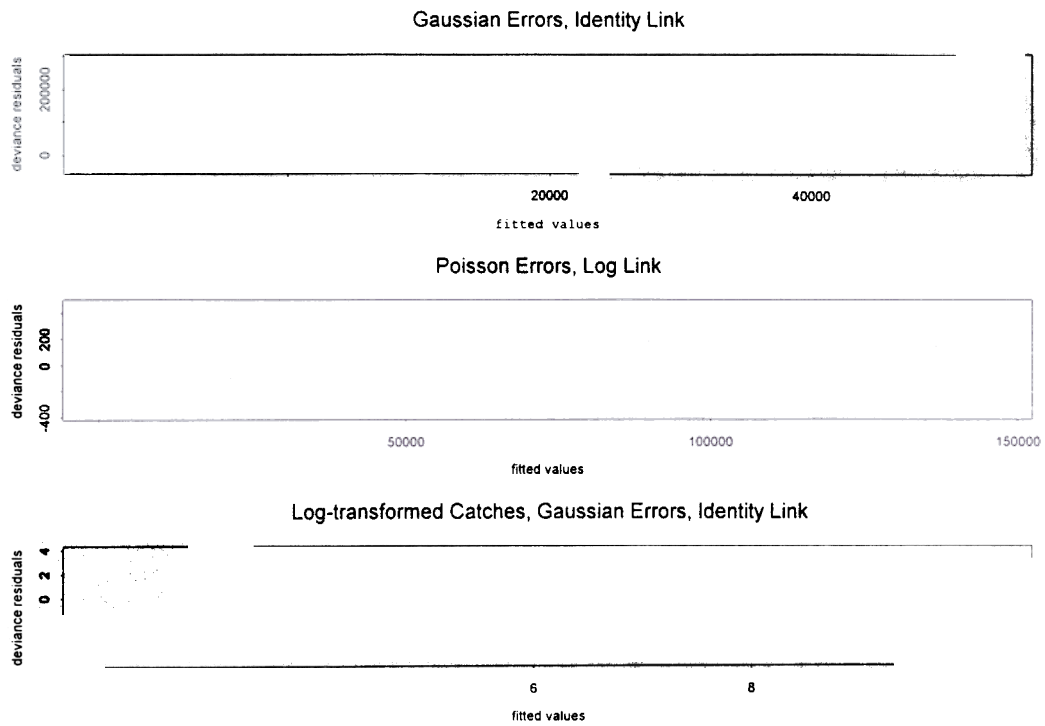| | Current | Separate | Combined | Sun elevation | Vessel | ICA |
|---|---|---|---|---|---|---|
| Mean std dev | 0.53 | 0.76 | 0.22 | 0.22 | 0.26 | 0.18 |

**Acknowledgements**

**References**

Augustin, NH, Borchers, DL, Clarke, ED & Buckland, ST 1998 Spatio-temporal modelling for the annual egg production method of stock assessment using generalized additive models. *Canadian Journal of Fisheries & Aquatic Sciences,* **55,** 2608-262 1.

Borchers, DL, Buckland, ST, Priede, IG & Ahmadi, S 1997 Improving the precision of the daily egg production method using generalized additive models. *Canadian Journal of Fisheries & Aquatic Sciences, 54,* 2727-2742.

Burnham, KP and Anderson, DR 1998 Model Selection and Inference. Springer-Verlag, New York.

Davison, AC and Hinkley, DV 1997 Bootstrap Methods and Their Application. Cambridge University Press! Cambridge

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

Hastie, TJ and Tibshirani, RJ 1990 Generalized Additive Models. Chapman & Hall, London.

Heessen HJL, Dalskov, J and Cook, RM 1997 The International Bottom Trawl Survey in the North Sea, Skagerrak and Kattgat. *ICES CM 1997/Y:31*

ICES 2000 Report of the herring assessment working group for the area south of 62N. *ICES 2000/ACFM:10.*

Patterson, KR and Melvin, GD 1996 Integrated Catch at Age Analysis Version 1.2. *Scottish Fisheries Research Report No 38*

Swartzman, G, Huang, C and Kaluzny, S 1992 Spatial analysis of Bering Sea groundfish survey data using generalized additive models. *Canadian Journal of Fisheries and Aquatic Sciences, 49,* 1366- 1378.

Swartzman, G, Stuetzle, W, Kulman, K and Powojowski, M 1994 Relating the distribution of pollock schools in the Bering Sea to environmental factors. *ICES Journal of Marine Science,* **51,** 48 1-492.

Swartzman, G, Silverman, E and Williamson, N 1995 Relating trends in Walleye Pollock *(Theraga chalcogrummu)* abundance in the Bering Sea to environmental factors. *Canudiun Journal of Fisheries and Aquatic Sciences, 52,* 369- 380.

Swartzman, G 1997 Analysis of the summer distribution of fish schools in the Pacific Eastern Boundary Current. *ICES Journal of Marine Science, 54, 40 l-l 16.*

Venables, WN and Ripley,- BD 1997 Modern Applied Statistics with S-PLUS (2nd Edn). Springer-Vcrlng. New York.
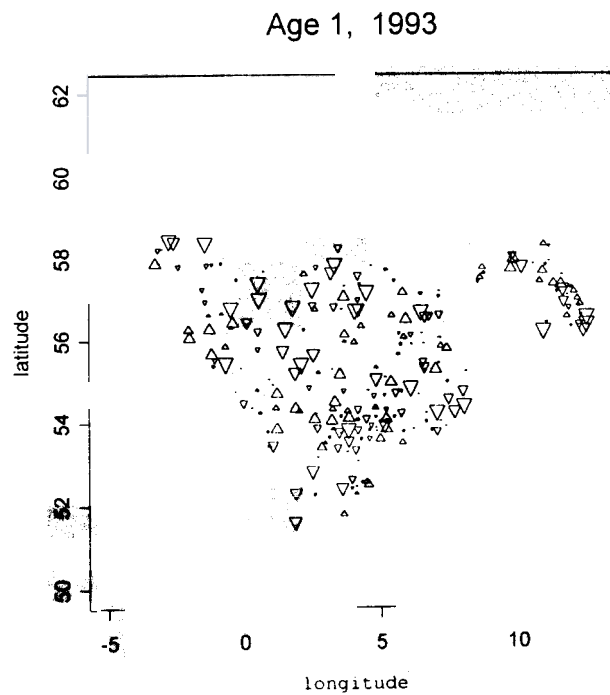
# 1993



**Figure**   The locations of the trawls in the 1993 survey.



**Figure 2** Histograms of a) the original catch data for 1993 age 1 herring, b) the log-transformed catches with a threshhold of 0.1 added so that zeros can be included in the plot and c) the log-transformed positive catches with a Gaussian distribution with the same mean and variance superimposed.
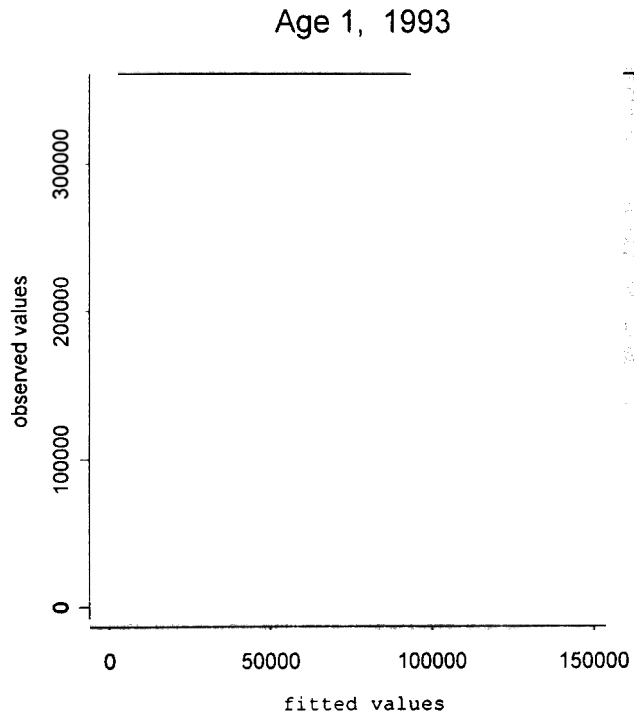
## Gaussian Errors, Identity Link

deviance residuals

20000          40000

fitted values

## Poisson Errors, Log Link

deviance residuals

50000          100000          150000

fitted values

## Log-transformed Catches, Gaussian Errors, Identity Link

deviance residuals

6               8

fitted values

**Figure 3** Residual plots for age 1 herring in 1993 using a) Gaussian errors with an identity link, b) Poisson errors with a log link and c) log-transformed data with Gaussian errors and an identity link.

## Age 1, 1993

latitude

-5          0          5          10

longitude

**Figure 4** Spatial residuals for positive catches of age 1 herring in 1993 using log-transformed data with Gaussian errors and an identity link. Upwards pointing triangles indicate positive residuals, downwards indicate negative residuals, the magnitude of the residual increasing with triangle size.

## Age 1, 1993



**Figure 5** Observed catches against predicted catches for age   herring in 1993 using the two-stage model.



**Figure 6** Fitted smooths for positive catches of age   herring in 1993 using log-transformed data with Gaussian errors and an identity link.

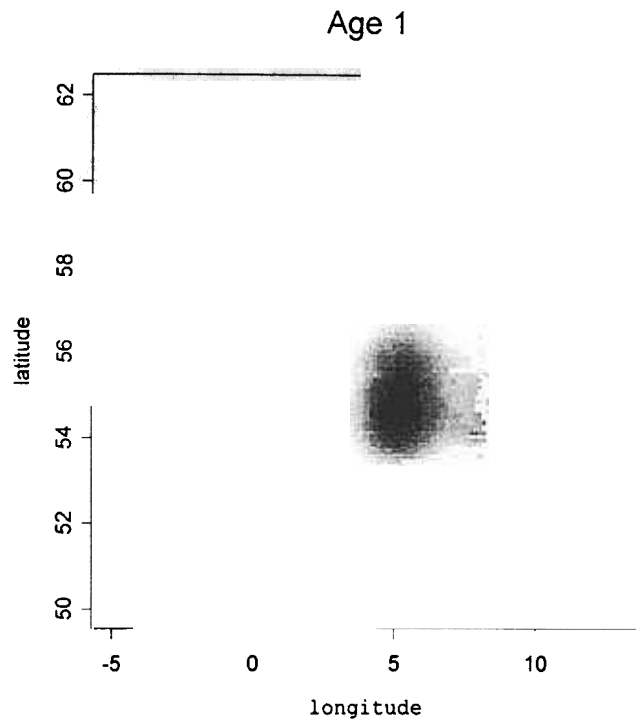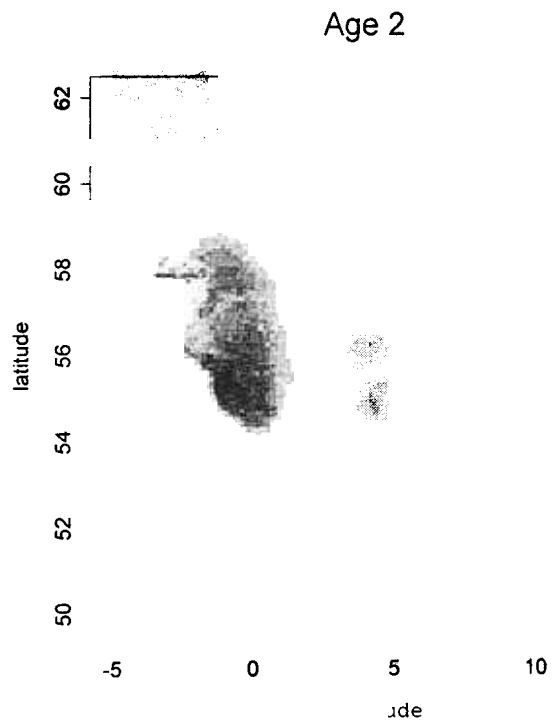**Figure 7a** Vessel effect for the presence of age 1 herring over 1983-1997.



**Figure 7b** Vessel effect for positive catches of age 1 herring over 1983-1997.

## Age 1



**Figure 8a** The generic spatial distribution of for age 1 herring.

## Age 2



**Figure 8b** The generic spatial distribution of for age 2 herring.

Δg

igu  8c Th        spati   istribution of for age  herrir

Ag

⊢

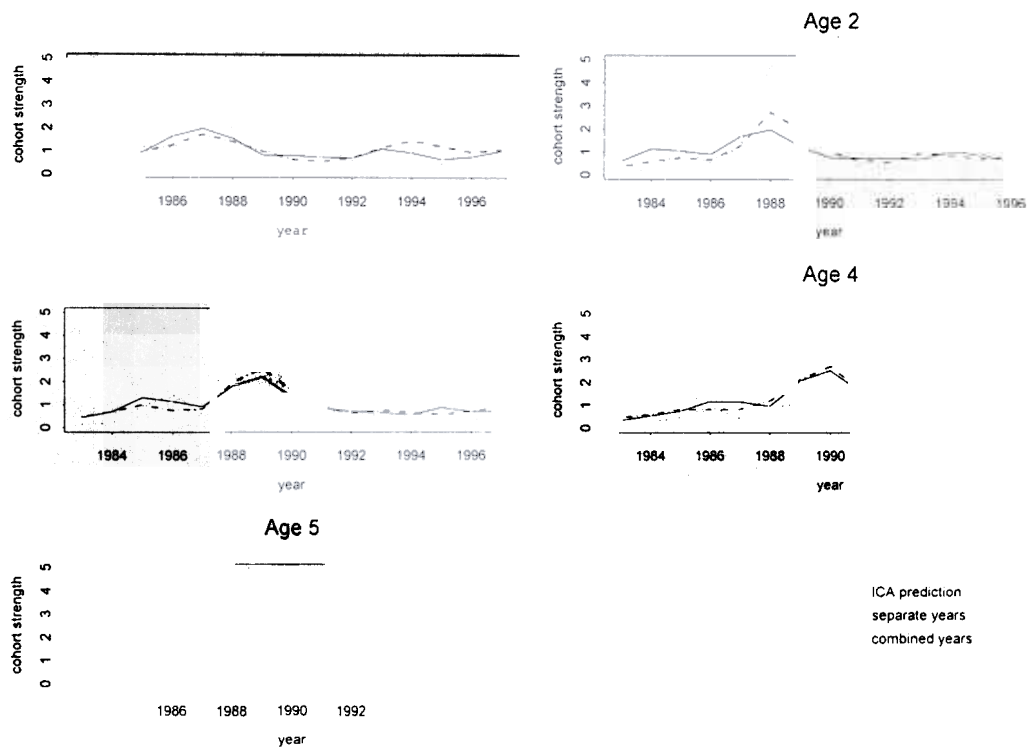gu        eric           ion           ir

# Age 5



**Figure 8e** The generic spatial distribution of for age 5 herring.

**Figure 9** Time series of cohort strength calculated using the current method and the best GAM-based model for each separate year, compared with the ICA predictions.



**Figure 10** Time series of abundance estimates calculated using the best model for each separate year and the model for all years combined, compared with the ICA predictions.
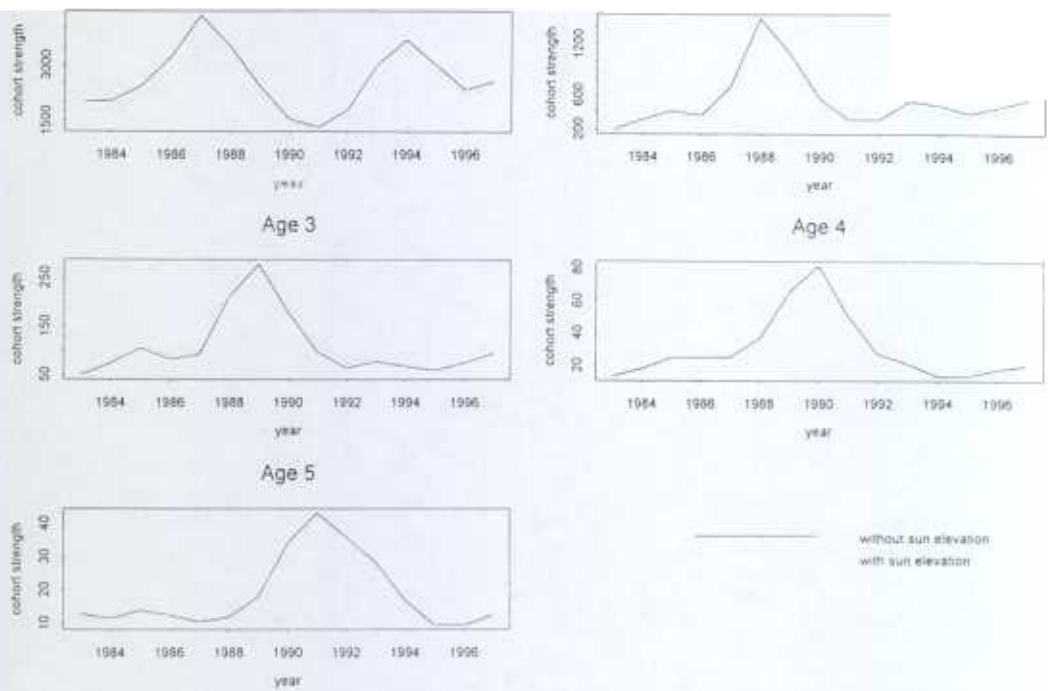
Age 3

Age 4

Age 5

**Figure 11a** Time series of abundance calculated using the model for all years combined, with and without sun elevation.

Age 1

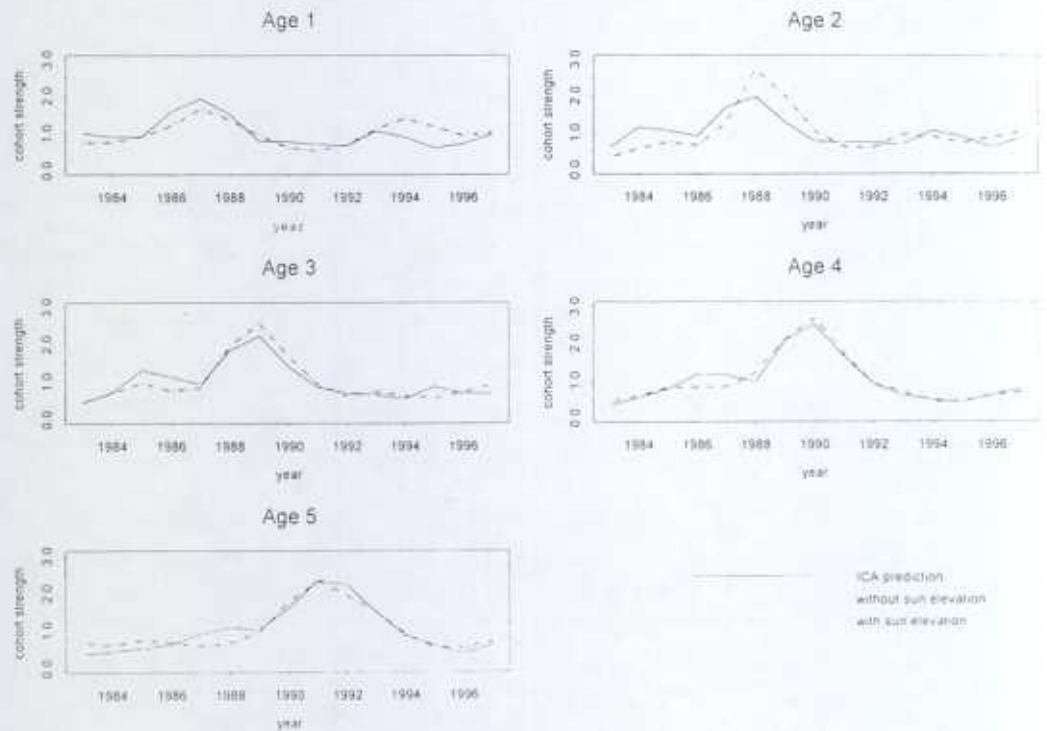Age 2

Age 3

Age 4

Age 5

**Figure 11b** Time series of cohort strength calculated using the model for all years combined, with and without sun elevation, compared with the ICA predictions.

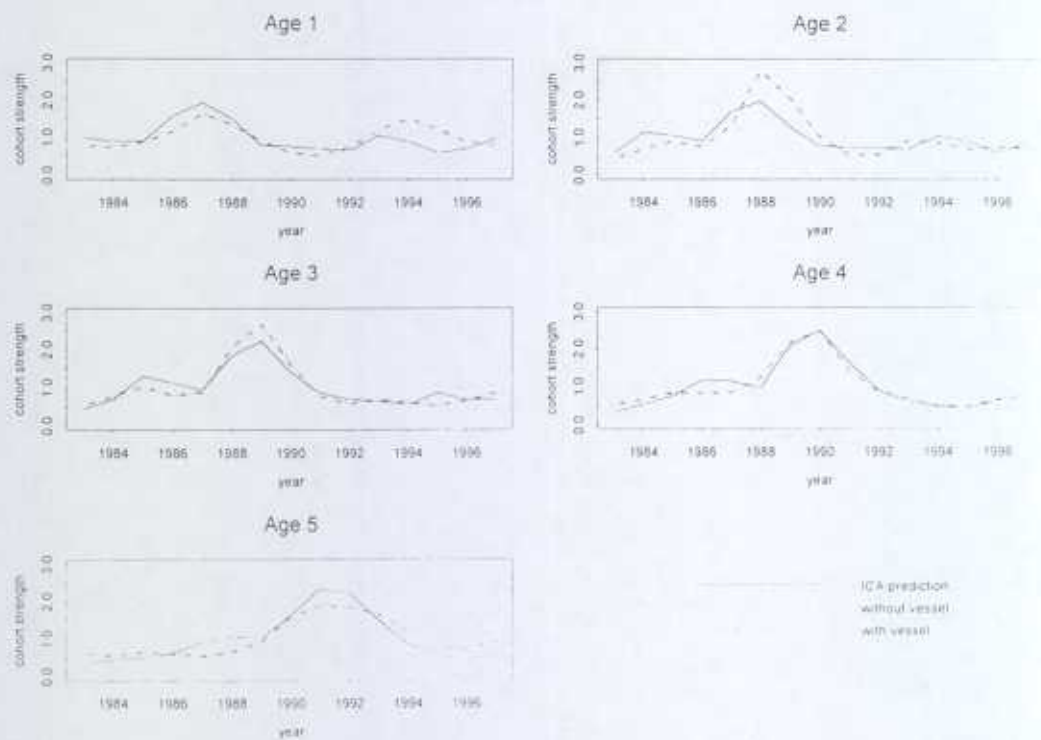without correcting for vessel effects.

## Age 1



## Age 2



## Age 3



## Age 4



## Age 5



ICA prediction
without vessel
with vessel

**Figure 12b** Time series of cohort strength calculated using the model for all years combined, with and without correcting for vessel effects, compared with the ICA predictions.