

# Analysis of groundfish survey abundance data: combining the GLM and delta approaches

Gunnar Stefánsson



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

This paper describes a method for the analysis of groundfish survey data by incorporating zero and non-zero values into a single model. This is done by using a model which modifies the delta-distribution approach to fit into the GLM framework and uses maximum likelihood to estimate parameters. No prior assumptions of homogeneity are used for the structure of the zero or non-zero values. The method is primarily applicable to fixed-station designs, although extensions to other designs are possible. The maximum likelihood estimation reduces to fitting a GLM to 0/1 values and another GLM to the positive abundance values. The new model is tested on Icelandic groundfish survey data. It is seen that the model can be used for evaluating the effect of different factors on catch rates as well as estimating abundance indices. Results from different models are compared on the basis of tuned VPA runs.

© 1996 International Council for the Exploration of the Sea

Key words: generalized linear models, marine surveys.

Received 7 September 1994; accepted 1 June 1995.

Gunnar Stefánsson: Marine Research Institute, P.O. Box 1390, Skúlagata 4, 121 Reykjavík, Iceland.

## Introduction

Groundfish surveys are commonly used for the purpose of obtaining an average catch per tow, to be used as an indicator of stock abundance. Data from individual tows have long been known to be notoriously variable (Thompson, 1928) and, hence, there is quite a long history of methods relating to the design and analysis of survey data. The emphasis in this paper will be on analysis techniques, and the method to be developed will be applicable in principle to most designs.

Several entirely different approaches exist for the analysis of groundfish survey data. Some of these methods are intricately linked to the design of the survey and can be classified according to whether the design is based on randomised or fixed stations. The methods can also be classified according to assumptions of the spatial distribution of the species and those made on the probability distribution of the measurements. Similar techniques have been developed and tested for acoustic survey data.

Most methods for the analysis either assume a homogeneous population within some regions (or strata) or estimate a single average within the stratum. Thus, within each stratum, the assumption is that all the measurements are of the same average population mean.

When stations are randomized every year, this assumption is true to some extent, although it usually wastes information by ignoring station location and in no way acknowledges the fact that there is always an underlying spatial pattern to the fish density, often with some year-to-year consistency. On the other hand, a randomized design should be set up in such a fashion as to incorporate the spatial information at the design stage, during the definition of strata.

The analysis then boils down to evaluating an average within each stratum and integrating these averages to obtain a stock index for the whole region.

Probably the single most common method for the analysis is the stratified analysis of Cochran (1977), commonly used for the analysis of data where stations have been allocated using a stratified random design (Smith, 1990). Alternatives include the so-called delta-distribution (Aitchison, 1955; Pennington, 1983) where zero values are treated separately and positive values are assumed to follow a lognormal distribution. The Adès distribution (Perry and Taylor, 1985) can be considered a relative of the delta distribution. As before, no spatial pattern is allowed within the strata. The delta-distribution would be better named the delta-lognormal distribution as it is perfectly feasible to use a similar delta-gamma distribution (Steinarsson and Stefánsson, 1986).

Entirely different approaches have also been tried for this type of general data, i.e. data on the amount in numbers or weight of fish caught per tow by commercial or research vessels. These include the use of models which assume linearity on a logarithmic scale (Gavaris, 1980; Myers and Pepin, 1986; Large, 1992), kriging (Petigas, 1993), linear models on a log-scale with spatially correlated errors (Polacheck and Volstad, 1993), generalized linear models (Smith *et al.*, 1991), and generalized additive models (Swartzman *et al.*, 1992). Similarly, several approaches have been attempted for the analysis of acoustic data, as detailed in Foote and Stefánsson (1993). In all these approaches, the underlying spatial distribution is or can be explicitly modelled and these methods for analysis have been used for data obtained from various survey designs. However, these particular methods tend to have problems with zero values. In particular, when the data from each tow are split into age groups, a large number of zero values can occur. Many of these exist simply because the tows are taken far away from the potential location of this particular age group. Other zero values may be important indicators of small stock size. Thus, one should consider models where these two types of zero values automatically influence the biomass indices in the right way. It is usually not possible to limit exactly the area of interest and this may have severe effects on the stock estimates for some procedures of analysis. Further, the fact that the data are best analysed in an age-disaggregated fashion compounds problems inherent in log-transforms (Myers and Pepin, 1990), since disaggregation will be likely to lead to many low abundance values, if there are several age groups in the stock of interest.

This paper develops a maximum likelihood method where an explicit formula is written down for the probability distribution of the catch at each station. This distribution will incorporate all the considerations mentioned above. The resulting model allows formal testing of what factors influence survey catches as well as the computation of abundance indices.

Data on haddock from an annual Icelandic groundfish survey (Pálsson *et al.*, 1989) will be used to exemplify the method.

## Data sets

The Icelandic groundfish survey has been conducted annually in March since 1985. Station locations are fixed in principle although minor variations can occur and, although the survey design initially included some 600 stations, only 488 have been taken in every one of the 10 survey years giving a total of 4880 observations for the current analysis. A detailed description of the survey is given in Pálsson *et al.* (1989). At each station, most fish are measured for length and at a number of stations samples are taken for ageing.

The method to be developed used as the basic datum the number of fish of a given age group at a given station. For this purpose, age samples within a stratum were pooled in order to obtain an age-length key, which was used to age-disaggregate the total length distribution at each station. The basic strata considered are given in Figure 1. The data on Icelandic haddock were obtained by smoothing the age-length key for each of the 10 areas and applying the resulting key to the length distribution at each station to obtain numbers at age by station and year.

Abundance estimates are taken from Anon. (1994). Since survey data were used to tune the abundance estimates, the years 1992–1994 were omitted and only the years 1985–1991 were used for comparing indices with VPA-based abundance estimates.

## The distribution

### Overall distribution of numbers per tow

Typical histograms of total or age-disaggregated catches exhibit considerable skewness and may have a spike at zero. As seen from the examples of basic data summaries given in Figures 2, 3, there may be a large number of zero values and a heavy tail. Hence, it is often advantageous to present data on a logarithmic scale. It is also often illustrative to consider histograms using  $\log(x+0.001)$  or a similar transformation, to see whether or not the zero values form a natural part of the whole. The choice of a low additive constant (0.001) is deliberate in this instance since it will isolate the zero points much more clearly than the use of a value closer to the smallest data value. In particular, in Figure 2, the zero values do not seem to be a natural part of the same continuous distribution as the one producing the positive values.

When zero values are eliminated, it is seen that the data may be close to lognormal, which again implies that a lognormal or gamma density may be appropriate for positive values, or possibly that a negative binomial may suffice for the entire data set. The negative binomial distribution, however, has a built-in linkage between the probability of zero and the mean of the positive values. This linkage will not hold when the area under consideration is changed to include more or fewer stations where the age group does not appear. Thus, the negative binomial distribution will not usually be applicable unless considerable attention is given to how zero tows are included in the analysis. The negative binomial distribution is a discrete distribution and thus might be believed to be applicable to count data as obtained in surveys (Taylor, 1953). However, the data are usually first scaled to tow duration or length and then disaggregated using proportions at age, and this will immediately lead to non-integer data.

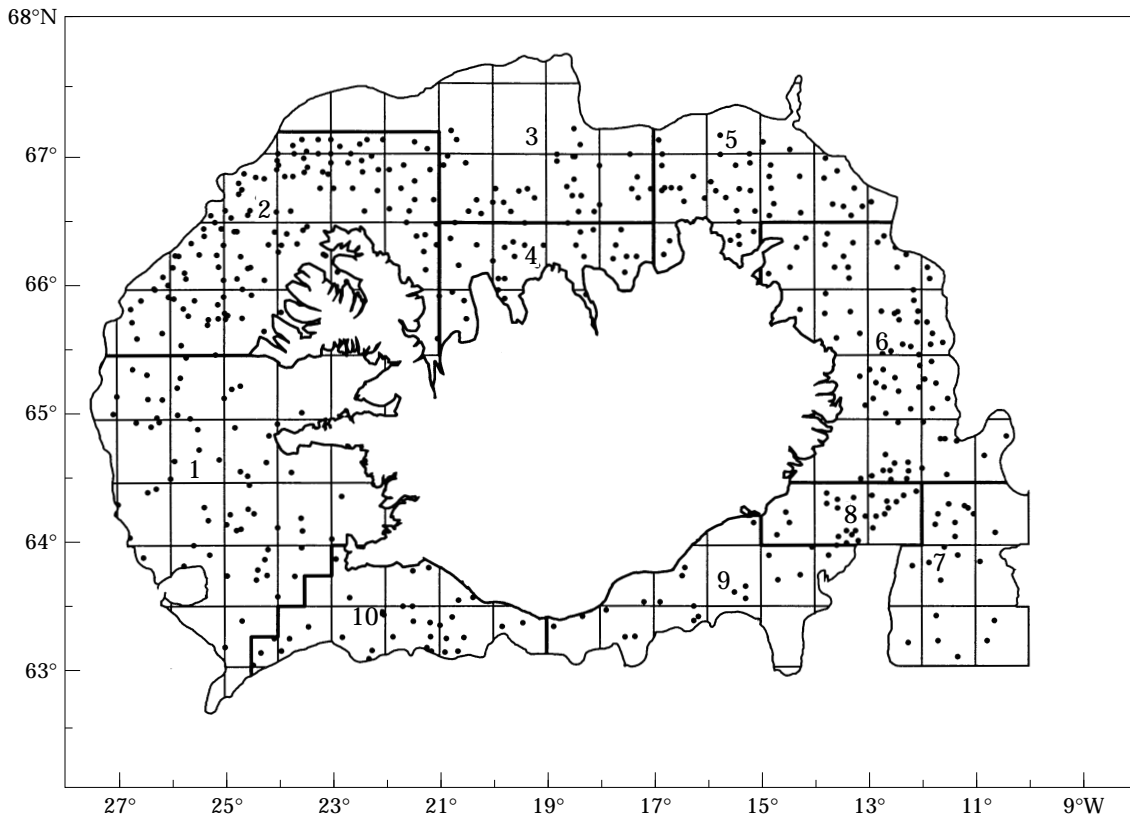


Figure 1. Stations, regions, and statistical squares used in the analysis.

The results of Myers and Pepin (1990) suggested that the use of the gamma density is preferable to the use of a lognormal density for fisheries data, although this seems to apply mainly when there is a considerable probability of small observations, not dealt with otherwise (Pennington, 1991) and, in other instances, the gain is minor (Firth, 1988). Although other members of the exponential family could be used, the gamma density is what will be used here when the positive values are under consideration, and it can be seen in Figures 2, 3 that the mean-variance relationship in the data seems to support the use of this density.

When a small year class appears, its distribution may change from the average in a number of ways. The density may stay constant at many points but the extent of the spatial distribution may diminish. This density change could thus result in the positive part of the histogram having the same mean but the number of zero values would increase. In the exact opposite case, the spatial distribution may stay the same but the density may go down at each point, though never to zero. These different types of changes have been investigated e.g. by Myers and Stokes (1989). A good mathematical model for the analysis of groundfish survey data should be able to accommodate these different conceptual models.

To account for the above considerations, the number of fish caught at a station,  $s$ , in year  $t$  may be taken to follow a distribution with a discrete probability of obtaining a non-zero count zero and some continuous density for the positive values. Thus, the cumulative distribution function (c.d.f.) of the abundance at the particular station becomes:

$$F_{st}(\omega) = P[Y_{st} \leq \omega] = (1 - p_{st}) + p_{st}G_{st}(\omega)$$

where  $G_{st}$  is a continuous c.d.f. describing the distribution of positive values. This is an extension of the general approach of Aitchison (1955), but the current framework will not assume that the parameters are constant from one station to the next.

When  $p_{st}$  is taken to be a constant within a stratum and  $G_{st}$  is a fixed lognormal distribution throughout the stratum, this is the usual delta-lognormal model (Pennington, 1983). If  $p_{st}$  is taken as the constant one and  $G_{st}$  is the negative binomial, we obtain another well known approach (Myers and Pepin, 1986). If zero values are omitted so that  $p_{st}$  is set to 1 and  $G_{st}$  is taken to be gamma density with a parameterized mean, this reduces to a generalized linear model (GLM).

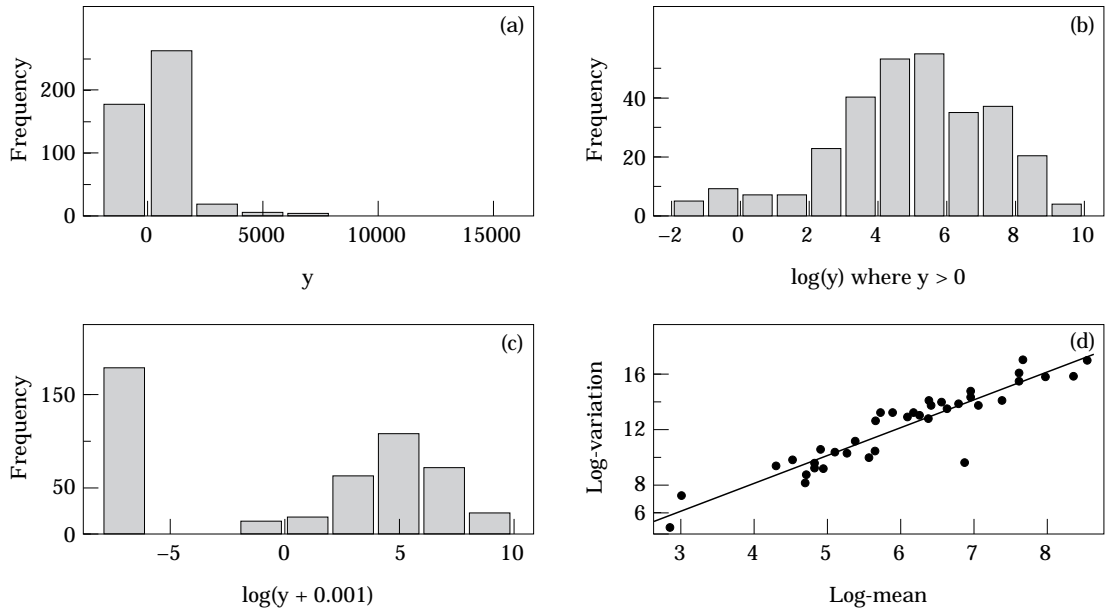


Figure 2. Data summaries of 2-group haddock abundance ( $y$ ) by tow from Icelandic groundfish survey. (a) Histogram of basic data values,  $y$ , (b) histogram of logged data values,  $\log(y)$  for  $y > 0$ , (c) histogram of logged slightly shifted data values,  $\log(y + 0.001)$ , (d) scatterplot of log-variance vs. log-mean for each statistical square with regression line.  $\text{Log}(\text{var}) = 0.16 + \log(\text{mean}) \times 2.00$ .

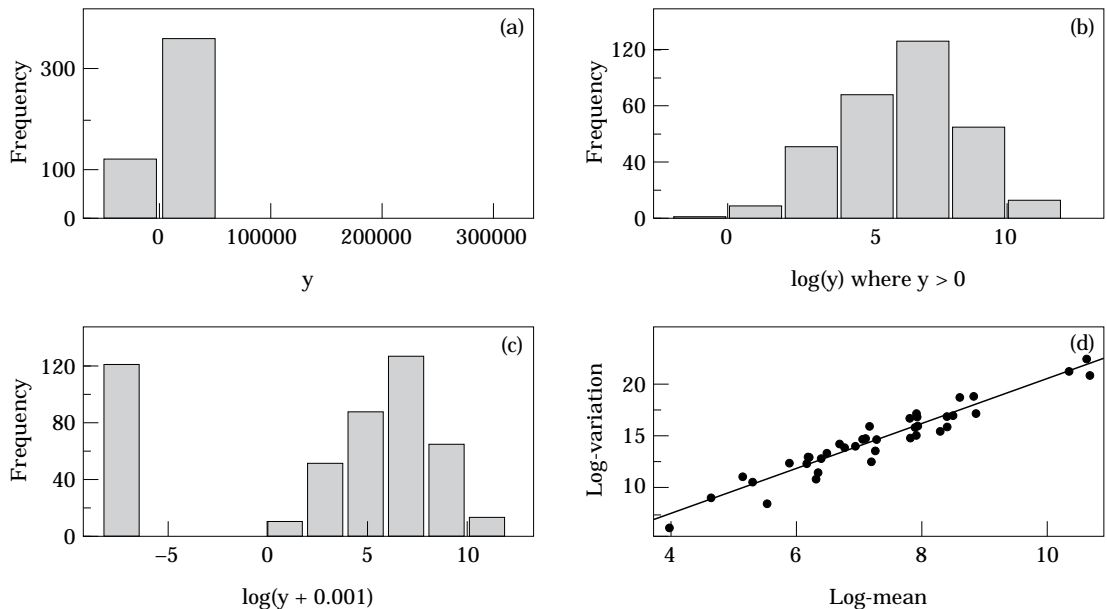


Figure 3. Data summaries of 3-group haddock abundance ( $y$ ) by tow from Icelandic groundfish survey. (a) Histogram of basic data values,  $y$ , (b) histogram of logged data values,  $\log(y)$  for  $y > 0$ , (c) histogram of logged slightly shifted data values,  $\log(y + 0.001)$ , (d) scatterplot of log-variance vs. log-mean for each statistical square regression line.  $\text{Log}(\text{var}) = -1.60 + \log(\text{mean}) \times 2.23$ .

From now on, a gamma density will be assumed for the positive values. The usual form of the gamma density function is given by:

$$\frac{y^{\alpha-1} e^{-y/\beta}}{\Gamma(\alpha)\beta^\alpha}$$

but, within GLMs, this is usually rewritten in terms of the mean,  $\mu=\alpha\beta$ , and one other parameter (the *shape* parameter),  $r$ , so that:

$$f(y) = \frac{y^{r-1} e^{-ry/\mu}}{\Gamma(r) \left[ \frac{\mu}{r} \right]^r}$$

The likelihood corresponding to the above c.d.f. is given by

$$L = \prod_{s,t:y_{st}=0} (1-p_{st}) \prod_{s,t:y_{st}>0} p_{st} \frac{y_{st}^{r-1} e^{-ry_{st}/\mu}}{\Gamma(r) \left[ \frac{\mu}{r} \right]^r}$$

Denote by  $n_{st}$  the number of repetitions of station  $s$  in year  $t$ , and by  $r_{st}$  the number of positive values obtained at this station. The above likelihood function can then be written in the form:

$$L = \prod_{s,t} (1-p_{st})^{n_{st}-r_{st}} p_{st}^{r_{st}} \prod_{s,t:y_{st}>0} \frac{y_{st}^{r-1} e^{-ry_{st}/\mu}}{\Gamma(r) \left[ \frac{\mu}{r} \right]^r}$$

It should be noted that, in almost all surveys,  $n_{st}=1$  and  $r_{st}$  is either 0 or 1. When repeated tows are performed at each station, the above formula must be understood to incorporate each station only once in the left part, but every positive number occurs once in the second product.

In the above formula there are two distinct components, the probability of a non-zero value and the distribution of the non-zero values. These can therefore be modelled and fitted separately to obtain first a fitted probability of non-zero tows and then the expected number of fish, given that some were caught. The fitted (unconditional) mean value at each station is then given by  $p_{st}\mu_{st}$ . Thus, the proposed model consists of two generalized linear models using a Bernoulli and a gamma distribution, respectively.

A model very similar to this one has been used for meteorological applications (Coe and Stern, 1982) and for data on consumption by cod (Waiwood *et al.*, 1991), although there are differences in model detail, particularly in the following model for the proportion. One virtue of the current approach in terms of fisheries applications is that it will allow formal testing of strata adequacy and several related issues.

### Modelling the probability

The probability of a non-empty tow can be modelled quite generally using generalized linear models (McCullagh and Nelder, 1989), as follows. The data is first recoded so that for each tow the value 0 is recorded if no fish are caught and the value 1 is recorded for non-zero tows to obtain Bernoulli-type 0/1-measurements. The usual model for probabilities is via the logit function, so that if the probability of a non-zero value is thought to depend on the latitude,  $h$ , then it would be appropriate to model the existence of fish in the trawl as a Bernoulli random variable with probability  $p$  given by:

$$\log(p/(1-p)) = \alpha + \gamma h$$

or, equivalently

$$p = \frac{1}{1 + e^{-(\alpha + \gamma h)}}$$

The formal statistical model is now to assume that the 0/1-values are independent results from measurements of a Bernoulli random variable with the probability,  $p$ , of ‘‘success’’, as given above. This model is in the class of generalized linear models which are available in some statistical packages such as Splus (Becker *et al.*, 1988; Chambers and Hastie, 1991).

In this setting, there is no particular reason to limit the linear predictor to specific variables and functions. Rather, the model should be thought of in the same light as ordinary regression models where parameters are tested for usefulness.

### Non-empty tows

In this section, only the non-zero tows are considered, so this analysis is conditional on the appearance of fish in the trawl. The basic model that will be used is the generalized linear model (GLM), where the number of fish is related to other measured variables through distributional assumptions.

In the gamma model, the variance is given by  $\sigma^2 = \mu^2/r$ , i.e. the variance is proportional to the square of the mean response (e.g. McCullagh and Nelder, 1989). It follows that if a gamma density is assumed, a regression slope based on a log-log plot of the variance vs. average abundance within predefined small (homogeneous) cells should be close to 2. In general, to determine the appropriate model, a base analysis of the raw data is required. An important item in such a base analysis is the above log-log plot, as exemplified in Figures 2, 3.

It is worth noting that the *shape* parameter  $r$  is related to the CV of the measurements through  $(CV)^2 = 1/r$ . This inverse is commonly called the *dispersion parameter* for the gamma family. As for the zero values, the GLM

approach relates the expected number of fish ( $\mu_{st}$ ) in a given year and location to other measurements through a link function, usually taken as log-linear in independent variables such as the length of the tow, a location effect and a year effect:

$$\ln \mu_{st} = \alpha + \beta l + \delta_s + \zeta_t$$

Naturally, factors and continuous variables can be used in an arbitrary mix in the two GLMs considered here, as in any GLMs.

### Combining the analyses

The resulting overall model for the distribution will be referred to as the Delta-Gamma ( $\Delta - \Gamma$ ) model.

Having obtained fitted values for the probability,  $p$ , of a non-zero tow, and for the expected number,  $\mu$ , conditional on it being positive, the predicted unconditional number of fish is given by  $p\mu$  and it should be noted that this quantity depends on whatever model was used for each of the two. For example, if the model for the proportion indicates that the probability of a non-zero tow differs between regions and the model for the mean of the non-zero values indicates that these depend on the length of the tow, then the overall mean,  $p\mu$ , depends on both the region and the length of the tow. The predicted tow content will be different for different lengths and regions.

### Model properties

It is possible to compute pointwise variances, i.e. the variance of the abundance at a station corresponding to fixed levels of the independent variables. The variance formula is obtained by appropriate (Riemann–Stieltjes) integration with respect to the cumulative distribution function,  $F(x)$  which has a point mass of  $1 - p$  at  $x=0$  and is  $p$  times a gamma c.d.f. at positive  $x$  values. The resulting variance is given by

$$\text{Var}(X) = p\sigma^2 + \mu^2 p(1 - p) = \mu^2 [p(1 + 1/r) - p^2]$$

and this is seen to be analogous to the results given in Aitchison (1955), Pennington (1983) and Smith (1988).

It follows that these pointwise variances can be estimated using the corresponding model estimates. It is seen that the overall variance is related to the overall mean, as in the simple gamma model, but with a different constant of proportionality, which depends on the proportion of non-zero tows.

It should be noted that this variance expression goes to a binomial variance as  $r$  goes to infinity and to a gamma variance as  $p$  goes to 1. When  $r > 1$ , the maximum (over  $p$ ) of this variance is obtained at

$$p = \frac{1}{2}(1 + 1/r)$$

where the variance becomes

$$\text{Var}(X) = \mu^2 [1 + 1/r]^2 / 4$$

When  $r < 1$  the maximum over  $p$  is obtained at  $p = 1$ , with a corresponding upper bound on the variance given by the gamma variance,

$$\text{Var}(X) = \mu^2 / r$$

These bounds are not sufficient to bound the coefficient of variation which is now defined by  $\sigma/p\mu$  where  $\sigma$  is the square root of the above variance. In fact, this CV can be made arbitrarily large by decreasing the proportion of non-zero tows.

Ideally, additivity should hold, since a property of towing is that, for a minor change in tow duration, a linear change in the abundance in the tow would be expected as a function of tow duration and the same overall distribution would be expected to hold for a long tow as for a short tow. Unfortunately, this is not the case. This can be seen by considering the characteristic function of the sum of two  $\Delta - \Gamma$  random variables. The characteristic function for the gamma density is given by

$$\varphi(t) = E[e^{itX}] = \int_0^{\infty} e^{itx} \frac{x^{r-1} e^{-rx/\mu}}{\Gamma(r)(\mu/r)^r} = \frac{1}{(1 - it\mu/r)^r}$$

and it follows that the characteristic function of the  $\Delta - \Gamma$  distribution is given by

$$\varphi^{GB}(t) = (1 - p) + \frac{p}{(1 - it\mu/r)^r}$$

The characteristic function of a sum of two independent random variables is the product of the two characteristic functions and it is therefore clear that the sum of two  $\Delta - \Gamma$  variables will not in general be a new  $\Delta - \Gamma$  variable. This is also clear if one considers the nature of the measurements themselves: if two identical stations are aggregated, then the result will be zero with probability  $(1 - p)^2$ , drawn from a single gamma density with probability  $2p(1 - p)$ , and it is, with probability  $p^2$ , drawn from a gamma density corresponding to the aggregate of two non-empty tows. It follows that the probability distribution of the sum of two  $\Delta - \Gamma$  variables is quite complex and analytical results are non-trivial to obtain. This conclusion leaves something to be desired since it would be useful to have an additive property for the distribution under consideration, particularly when considering varying tow duration. Investigating alternative c.d.f.s with such a property would seem to be a useful area for future work.

Table 1. Analysis of deviance table for different Bernoulli-based (Delta) generalized linear models fitted to presence/absence of 2-group haddock by station in Icelandic groundfish survey, 1985–1994. The models are fitted sequentially and the columns give the residual degrees of freedom for each model, the residual deviance, the degrees of freedom corresponding to the additional term, the resulting change in deviance, and the p-value when a  $\chi^2$ -test is used to test for significance.

Model terms	Residual d.f.	Residual deviance	Test d.f.	Change in deviance	p
1	4879	6528			
Year	4870	6432	9	96	<0.001
Year+ns	4869	6096	1	335	<0.001
Year+reg	4861	5540	8	556	<0.001
Year+reg+depth	4852	4540	9	1001	<0.001
Year+square+depth	4757	3761	95	779	<0.001
Year+square	4766	4259	–9	–498	<0.001

## Details of fitting and interpretation. Numerical example

The data on age 2 and 3 haddock are summarized in Figures 2, 3. When fitting models to data it is of importance to decide appropriately which parameters are important and, in general, it will be of interest to test for effects such as diurnal differences in catch rates, depth, or temperature. In what follows, only depth and spatial factors will be considered in addition to the year effect.

Table 1 gives stepwise analysis of variance results for comparing a sequence of Bernoulli models fitted to the age 2 haddock data. All effects tested are entered as factors and the sequence of factors is in the direction of increasingly detailed splits of the oceanic area, from a simple north–south split (regions 2–7 vs. 1, 9, 10 in Figure 1) through the strata in Figure 1, to statistical squares.

As for other generalized linear models, the *deviance* is used in much the same way as the sum of squares is used in ordinary regression. In particular, the reduction in deviance is related to the usual concept of  $r^2$ . It is seen that the explained variation in Table 1 is not a very large fraction (40%) of the total deviance, indicating that the location of 2-group haddock can only somewhat poorly be explained by the model, but this is considered further below.

Since the data are 0/1 measurements, a  $\chi^2$ -statistic is used to test for significance. In Table 1, it is seen that the initial (null) deviance of 6528 is reduced by 96 to 6432 by using a model with only a year effect. Although this reduction in deviance is small in relation to the total deviance, it is considerable in comparison to the degrees of freedom (9) expended and, therefore, it is highly significant. It is seen that a considerable further reduction is obtained by accounting for the difference in the north and south regions, but a further similar reduction is obtained by going to a 10-region split. The 10 regions

clearly do not capture all the depth information, since depth is a highly significant addition to the model. In addition to this, the statistical squares seem to be a highly significant and important addition to the model. Notably, there may be a difference between “significant” and “important”. The “importance” of a variable can be determined e.g. in terms of the proportional reduction in deviance. A statistically significant variable may in some instances be better left out of the model, if the amount of variation explained by the variable is small in relation to the complexity that it adds.

These test results indicate that the finest possible spatial scale is needed. In fact, the greatest portion of the deviance is explained at the last step, going from 10 areas to the statistical squares. It is also seen (last line of Table 1) that it is not possible to drop the depth term from the final model, as this is still significant, even when the square effect has been entered.

The final model in Table 1 has a deviance of 4259 on 4766 d.f. which might possibly be taken to indicate an adequate fit to the data. However, as indicated in McCullagh and Nelder (1989, p. 119), tests for lack-of-fit will not be valid based on the current sparse data set (unless the data can be collapsed into a smaller frequency table for a given model containing only factors).

Table 2 gives test results for the gamma portion of the model for age 2 haddock. It is seen that the model can be used to reduce the deviance from 15141 to 8664 and that for this model a fine spatial scale is also needed. For these tests, an F-statistic has been used, since the gamma model contains an unknown scale parameter.

Table 3 and 4 give similar results for age 3 haddock. The same overall conclusions hold in that a considerable reduction in deviance can be obtained using the model.

Although the above indicates that the use of a parametric model can be used to considerable benefit in terms of reduction of variability ( $r^2$  in the range 0.4–0.5), it should be noted that there is still considerable residual variation in the data. Thus, the estimated dispersion

Table 2. Analysis of deviance table for different gamma-based generalized linear models fitted to the abundance of 2-group haddock. Data used consists of those stations with some catches of 2-group haddock in Icelandic groundfish survey, 1985–1994. The models are fitted sequentially and the columns give the residual degrees of freedom for each model, the residual deviance, the degrees of freedom corresponding to the additional term, the resulting change in deviance, and the p-value when an F-test is used to test for significance.

Terms	Residual d.f.	Residual deviance	d.f.	Change in deviance	F-value	p
1	2975	15 141				
Year	2966	13 263	9	1878	69.0	<0.001
Year+ns	2965	13 160	1	103	33.9	<0.001
Year+reg	2957	10 947	8	2213	91.4	<0.001
Year+reg+depth	2949	10 425	8	522	21.6	<0.001
Year+square+depth	2863	8664	86	1760	6.8	<0.001
Year+square	2871	8943	–8	–279	11.5	<0.001

Table 3. Analysis of deviance table for different Bernoulli-based (Delta) generalized linear models fitted to presence/absence of 3-group haddock by station in Icelandic groundfish survey, 1985–1994. The models are fitted sequentially and the columns give the residual degrees of freedom for each model, the residual deviance, the degrees of freedom corresponding to the additional term, the resulting change in deviance, and the p-value when a  $\chi^2$ -test is used to test for significance.

Terms	Residual d.f.	Residual deviance	d.f.	Change in deviance	p
1	4879	5965			
Year	4870	5896	9	70	<0.001
Year+ns	4869	5368	1	527	<0.001
Year+reg	4861	4847	8	522	<0.001
Year+reg+depth	4852	4031	9	816	<0.001
Year+square+depth	4757	3116	95	915	<0.001
Year+square	4766	3542	–9	–425	<0.001

Table 4. Analysis of deviance table for different gamma-based generalized linear models fitted to the abundance of 3-group haddock. Data used consists of those stations with some catches of 3-group haddock in Icelandic groundfish survey, 1985–1994. The models are fitted sequentially and the columns give the residual degrees of freedom for each model, the residual deviance, the degrees of freedom corresponding to the additional term, the resulting change in deviance, and the p-value when an F-test is used to test for significance.

Terms	Residual d.f.	Residual deviance	d.f.	Change in deviance	F-value	p
1	3413	15 108				
Year	3404	12 746	9	2362	106	<0.001
Year+ns	3403	12 674	1	71	29	<0.001
Year+reg	3395	10 817	8	1857	94	<0.001
Year+reg+depth	3387	10 155	8	662	34	<0.001
Year+square+depth	3301	8144	86	2012	9	<0.001
Year+square	3309	8526	–8	–382	19	<0.001

parameters in the gamma models are 4.2 and 3.9 for age groups 2 and 3, respectively, indicating that the CV of the positive measurements is about 200% after correcting for the relevant factors and, for practical purposes,  $r$  can be taken to be about 0.25. In this context it is worth noting that computing the CV based on the raw positive

values yields some 470% in both cases. From the previous theory, it follows that the variance of the number of fish in a tow is  $\text{Var}(X) = \mu^2[p(1+1/r) - p^2]$  which is bounded by  $\text{Var}(X) = \mu^2/r$ . Thus, the CV of the catch will be no more than  $\sqrt{1/pr}$ . Taking the above value for the dispersion parameter, and considering a region where



there is more than 50% probability of catching haddock, the CV of the catch will be no more than 280%, but if haddock are always caught in a certain location, the CV will be no more than 200%.

It is also clear from the above examples how the current approach can be used to test or modify proposed stratification schemes. In particular, the deviance analyses will illustrate the importance of certain methods of stratification. Some measured variables may be important although they are not directly linked to the stratification method. For example, it is quite possible that depth, time of day, or other variables may be used to reduce the amount of variation in the data and thus potentially increase the usefulness of abundance indices, for example.

The results from the above analyses can be used alone to provide year effects (main effects) in the models indicated. Naturally, if main effects are to be extracted and used, special care needs to be taken if the models contain interaction terms with year. The common approach (Anon., 1992) of extracting the year effect alone from a log-linear model, for example, can be replaced by an integral of the fitted model over the entire region under consideration. This approach yields abundance indices which are equivalent to the year effects when the model contains no interaction terms. Within the present model, however, even when the year effect is only present as a main effect in each piece of the  $\Delta - \Gamma$  model, the overall model will contain a non-linear year effect. Thus, integration of the fitted response surface is required to obtain the overall abundance index.

Fitted values from the overall  $\Delta - \Gamma$  model are obtained from the two submodels. These values can be computed on a grid and then integrated. The present approach is simply to compute the average fitted value across stations within each statistical square and the overall integral is taken as the direct average over all squares. The resulting indices is analysed along with other indices of abundance in the following sections.

The binary nature of the Bernoulli part of the model leads to nontrivial problems in interpretation which are worthy of some note. In particular, when factors only are used in the model, it is possible to collapse the data into a contingency table which leads to an anova-table which contains different deviance values, although the test statistics are the same. To give an example, if only the effect of the north-south stratification is considered along with the year effect, than Table 5a,b is obtained based on Bernoulli (0/1) or binomial (collapsed) data, respectively. It is seen that in the binomial table a considerable proportion (99%) of the deviance is explained whereas only a minor proportion (10%) is explained in the Bernoulli table. It is also seen (by comparing the residual d.f. to the residual deviance) that the goodness of fit test results are considerably different. The reasons for these apparent discrepancies are of

Table 5. Comparison of analyses based on (a) the full data set (Bernoulli GLM), and (b) a collapsed set (binomial GLM). Icelandic groundfish survey. Age 3 haddock. Results are based on adding terms sequentially.

(a) Bernoulli model

	d.f.	Change in deviance	Residual d.f.	Residual deviance	p ( $\chi^2$ )
NULL	—	—	4879	5965.41	—
Year	9	69.77	4870	5895.65	<0.001
ns	1	527.24	4869	5368.41	<0.001

(b) Binomial model

	d.f.	Change in deviance	Residual d.f.	Residual deviance	p ( $\chi^2$ )
NULL	—	—	19	601.92	—
Year	9	69.77	10	532.14	<0.001
ns	1	527.24	9	4.90	<0.001

course that the Bernoulli model attempts to explain data on a much finer scale than does the binomial model.

## Other indices

A variety of different methods have been proposed to compute indices based on marine survey data. The simplest possible method is simply to compute the arithmetic mean (AM) over all the stations. Alternatives include a stratified mean (SM) using e.g. the 10 regions in Figure 1 and a geometric average (GM) based on

$$GM_i = e^{\frac{\sum \ln(y_{ni} + 1)}{s}} - 1.$$

Table 6a,b lists all these indices for age groups 2 and 3 of haddock, along with the  $\Delta - \Gamma$  index and VPA estimates of year class strength at ages 2 and 3, as given in Anon. (1994).

## Comparisons among indices and with VPA

The indices can be compared with each other and with population estimates using virtual population analysis (VPA) (Gulland, 1965; Anon., 1994) using correlation analysis and plots. The plots are given in Figures 4, 5 and  $r^2$ -values are given in Table 7. Also given in Table 7 is the result from predicting VPA values based on simple scaling of the indices. This is done by computing average catchability by dividing average VPA by the average index and multiplying the annual index values by catchability.

Table 6. Indices of abundance of 2- and 3-group Icelandic haddock, based on different methods of analysis by year and year class (Ycl). DG=delta-gamma, AM=arithmetic mean, SM=stratified mean, GM=geometric mean. VPA=virtual population analysis.

Indices of 2-group abundance										
Year	85	86	87	88	89	90	91	92	93	94
Ycl	83	84	85	86	87	88	89	90	91	92
DG	982	1821	3891	594	424	690	2390	4411	1073	1225
AM	604	1674	4816	624	357	615	2761	3820	566	889
SM	583	1671	4777	604	395	625	3092	4214	623	1053
GM	25	47	56	12	8	11	60	95	22	18
VPA	41	88	164	46	26	26	113	167	40	50

Indices of 3-group abundance										
Year	85	86	87	88	89	90	91	92	93	94
Ycl	82	83	84	85	86	87	88	89	90	91
DG	260	1580	2917	3508	717	518	836	2649	5320	802
AM	303	1067	3413	3575	732	504	718	2790	4455	622
SM	287	1013	3251	3363	835	560	742	3074	5046	684
GM	17	75	125	110	32	13	41	113	142	35
VPA	16	33	72	132	38	21	21	90	135	33

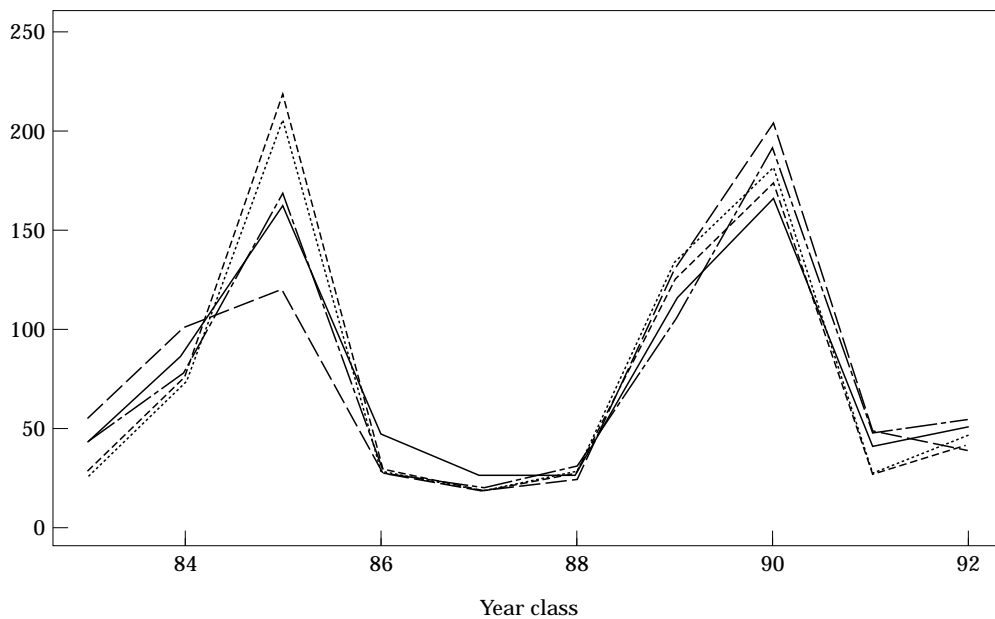


Figure 4. Time series of indices of 2-group haddock, scaled to VPA average. (—)=VPA; (· · ·)=GB; (· · ·)=SM; (---)=AM; (— —)=GM.

The various indices can be compared to results from tuned VPA estimates of stock size. In this context it must be noted that survey data play an important role in the estimate of stock size. To avoid possible confounding effects, comparisons with VPA are based only on data from the years 1985–1990.

The tables and figures show that there is no clear winner in the comparisons, although the  $\Delta - \Gamma$  index

does best in terms of predicting VPA values based on simple catchability scaling. Similarly, the  $\Delta - \Gamma$  index seems to provide very consistent indices for the two age groups. Although care should be taken not to over-interpret these results, it should be noted that these two properties are very important, since an index should at a minimum be internally consistent and it should ideally provide a proportional relationship with

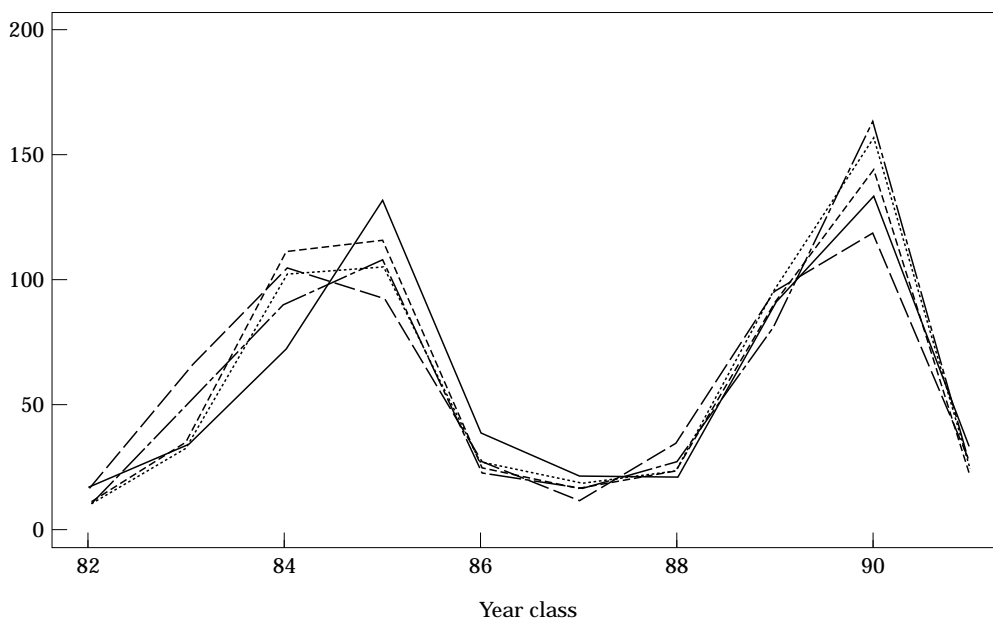


Figure 5. Time series of indices of 3-group haddock, scaled to VPA average. Key as for Figure 4.

Table 7. (a, b) Explained VPA variation ( $r^2$ ), as obtained from linear regression of VPA on the various indices using raw or log scale for 2- and 3-group haddock. Also given is the sum of squared errors (SSE) based on predicting VPA numbers from indices using a scale factor alone. (c) Explained variation when an index for age 2 is used to explain the index for age 3. DG=delta-gamma model indices, AM=arithmetic mean, SM=stratified mean, GM=geometric mean.

Index	(a) 2-group			(b) 3-group		
	$r^2$ log	$r^2$ raw	SSE	$r^2$ log	$r^2$ raw	SSE
DG	0.89	0.97	750	0.86	0.85	1487
AM	0.92	0.96	3573	0.92	0.83	2053
SM	0.91	0.96	3503	0.92	0.82	1906
GM	0.85	0.85	2101	0.75	0.63	3739

(c)  $r^2$  between log indices for ages 2-3

	DG3	AM3	SM3	GM3	VPA3
DG2	0.92	0.88	0.88	0.80	0.91
AM2	0.87	0.91	0.91	0.74	0.94
SM2	0.86	0.89	0.91	0.71	0.94
GM2	0.93	0.89	0.90	0.86	0.89
VPA2	0.87	0.90	0.92	0.72	1.00

abundance, rather than a non-linear one or one with an intercept.

As seen from Figures 4, 5, however, the results from the comparisons with VPA might change if another set

of years were to be used, but this is not possible within the present framework since the final VPA values need to be tuned to certain indices.

## Conclusions

The approach considered is based on a model which has considerable intuitive appeal in that it incorporates most of the concerns usually raised in the analysis of groundfish survey data. Results obtained are close to those obtained by other methods (which is in accordance with the results in Anon., 1992) but the *ad hoc* nature of many other methods is eliminated by using an explicit model for zero and non-zero values. The model provides an analysis technique where many problems usually associated with zero values are alleviated. This includes issues such as those involving the definition of an appropriate area for the analysis and those related to log-transforming values which can be arbitrarily close to zero. Furthermore, the approach can accommodate spatial and temporal variability in an explicit model.

Variance estimates for the resulting parameters are available but should be viewed with caution, since the real variances of interest are those related to prediction capabilities and the degrees of freedom vary depending on the inclusion of zero-catch tows. The actual variances of interest are probably better obtained by tuning VPAs with the indices, as in Anon. (1992).

This model has considerable potential for the general analysis of groundfish survey data, since it can incorporate several relevant properties of fish distributions, including changes in density and range. The usual

qualities of GLMs, specifically the potential for incorporating, estimating, and testing effects such as diurnal variations, are also available.

It must be noted, however, that here, as in Anon. (1992), the actual type of analysis considered does not seem to be of great consequence to the predictive power of abundance indices obtained, since even simple methods of analysis yield results fairly consistent (in a regression context) with VPA results for the data sets considered. This conclusion is, however, likely to be dependent on several factors such as the goodness of fit of the specific model and the number of stations in the survey.

One important item to note is that the issues raised in Pennington and Vølstad (1994) concerning the *intra* haul correlation and its effect on the various statistics are not considered at all in the present paper. The joint consideration of the coarse-scale spatial distribution modelled here and the finer-scale effects expressed in the *intra* haul correlation is of major interest and needs to be considered further.

Another important area of future work on the analysis of abundance data is in developing and testing scatterplot smoothers in generalized additive models (GAMs) as described by Hastie and Tibshirani (1990) where the abundance in a given location is described as a smooth function of independent variables such as location or depth.

## References

- Aitchison, J. 1955. On the distribution of a positive random variable having a discrete probability mass at the origin. *Journal of the American Statistical Association*, 50: 901–908.
- Anon., 1992. Report of the workshop on the analysis of survey data. ICES C.M. 1992/D:6.
- Anon., 1994. *Nytjastofnar sjávar og umhverfisþættir 1994. Aflahorfur fiskveiðiárið 1994/1995* [State of Marine Stocks and Environmental Conditions in Icelandic Waters 1994, Prospects for the Quota Year 1994/93]. Marine Research Institute, Reykjavik.
- Becker, R. A., Chambers, J. M., and Wilks, A. R. 1988. *The New S language. A programming environment for data analysis and graphics*. Wadsworth & Brooks/Cole Advanced Books & Software, Pacific Grove, California.
- Chambers, J. M. and Hastie, T. J. 1991. *Statistical models in S*. Wadsworth and Brooks, Pacific Grove, California. 608 pp.
- Cochran, W. G. 1977. *Sampling techniques*. 3rd edn. John Wiley & Sons, Inc., New York. 330 pp.
- Coe, R. and Stern, R. D. 1982. Fitting models to daily rainfall data. *Journal of Applied Meteorology*, 21: 1024–1031.
- Firth, D. 1988. Multiplicative errors: Log-normal or Gamma? *Journal of the Royal Society, B*, 50: 266–268.
- Foote, K. G. and Stefánsson, G. 1993. Definition of the problem of estimating fish abundance over an area from acoustic line transect measurements of density. *ICES Journal of Marine Science*, 50: 369–381.
- Gavaris, S. 1980. Use of a multiplicative model to estimate catch rate and effort from commercial data. *Canadian Journal of Fisheries and Aquatic Sciences*, 37: 2272–2275.
- Gulland, J. A. 1965. Estimation of mortality rates. Annex to Arctic Fisheries Working Group Report. ICES, C.M. Gadoid Fish Committee: 3. 9 pp.
- Hastie, T. and Tibshirani, R. 1990. *Generalized additive models*. Chapman and Hall, London. 335 pp.
- Large, P. A. 1992. Use of a multiplicative model to estimate relative abundance from commercial CPUE data. *ICES Journal of Marine Science*, 49: 253–261.
- McCullagh, P. and Nelder, J. A. 1989. *Generalized Linear Models*. Chapman and Hall, London. 511 pp.
- Myers, R. A. and Pepin, P. 1986. The estimation of population size from research surveys using regression models. ICES CM 1986/D:9.
- Myers, R. A. and Pepin, P. 1990. The robustness of lognormal based estimators of abundance. *Biometrics*, 46: 1185–1192.
- Pálsson, Ó. K., Jónsson, E., Schopka, S. A., Stefánsson, G., and Steinarsson, B. Æ. 1989. Icelandic groundfish survey data used to improve precision in stock assessments. *Journal of Northwest Atlantic Fisheries Science*, 9: 53–72.
- Pennington, M. 1983. Efficient estimators of abundance, for fish and plankton surveys. *Biometrics*, 39: 281–286.
- Pennington, M. 1991. On testing the robustness of lognormal-based estimators. *Biometrics*, 47: 1623–1624.
- Pennington, M. and Vølstad, J. H. (In press). Assessing the effect of intra-haul correlation and variable density on estimates of population characteristics from marine surveys. *Biometrics*.
- Perry, J. N. and Taylor, L. R. 1985. Adès: New ecological families of species-specific frequency distributions that describe repeated spatial samples with an intrinsic power-law variance-mean property. *Journal of Animal Ecology*, 54: 931–953.
- Petitgas, P. 1993. Geostatistics for fish stock assessments: a review and an acoustic application. *ICES Journal of Marine Science*, 50: 285–298.
- Polacheck, T. and Vølstad, J. H. 1993. Analysis of spatial variability of Georges Bank haddock (*Melanogrammus aeglefinus*) from trawl survey data using a linear regression model with spatial interaction. *ICES. Journal of Marine Science*, 50: 1–8.
- Smith, S. J. 1988. Evaluating the efficiency of the  $\Delta$ -distribution mean estimator. *Biometrics*, 44: 485–493.
- Smith, S. J. 1990. Use of statistical models for the estimation of abundance from groundfish survey data. *Canadian Journal of Fisheries and Aquatic Sciences*, 49: 1366–1378.
- Smith, S. J., Perry, R. I., and Fanning, L. P. 1991. Relations between water mass characteristics and estimates of fish population abundance from trawl surveys. *Environmental Monitoring and Assessment*, 17: 227–245.
- Steinarsson, B. Æ. and Stefánsson, G. 1986. Comparison of random and fixed trawl stations in Icelandic groundfish surveys and some computational considerations. ICES C.M. 1986/D:13.
- Swartzman, G. L., Huang, C. P., and Kaluzny, S. P. 1992. Analysis of Bering Sea groundfish survey data using generalized additive models. *Canadian Journal of Fisheries and Aquatic Science*, 49: 1366–1378.
- Taylor, C. C. 1953. Nature of variability in trawl catches. *Fishery Bulletin. U.S.*, 54: 145–166.
- Thompson, H. 1929. General features in the biology of the Haddock (*Gadus aeglefinus* L.) in Icelandic waters during 1903–1926. *Rapp. proc. verb. Reun. LVIII*: 3–73.
- Waiwood, K. G., Smith, S. J., and Petersen, M. R. 1991. Feeding of cod (*Gadus morhua*) at low temperatures. *Canadian Journal of Fisheries and Aquatic Sciences*, 48: 824–831.