

The efficient integration of abundance and demographic data

P. Besbeas,

University of Kent at Canterbury, UK

J.-D. Lebreton

Centre d'Ecologie Fonctionnelle et Evolutive, Montpellier, France

and B. J. T. Morgan

University of Kent at Canterbury, UK

[Received May 2001. Revised June 2002]

Summary. A drawback of a new method for integrating abundance and mark–recapture–recovery data is the need to combine likelihoods describing the different data sets. Often these likelihoods will be formed by using specialist computer programs, which is an obstacle to the joint analysis. This difficulty is easily circumvented by the use of a multivariate normal approximation. We show that it is only necessary to make the approximation for the parameters of interest in the joint analysis. The approximation is evaluated on data sets for two bird species and is shown to be efficient and accurate.

Keywords: Census data; Computational efficiency; Kalman filter; Multivariate normal approximation; Ring–recovery data; State space

1. Introduction

Animal population growth and change can be studied in two different ways: first, the population kinetics can be studied by successive population censuses, from which, for instance, estimates of growth rate can be deduced; secondly, demographic studies can provide estimates of demographic parameters, such as annual survival probabilities, which can in turn be used in population models to predict population trajectories. This is typically done by using Leslie matrices (Leslie, 1945; Caswell, 2000). Matching these two types of analyses is done in an *ad hoc* way, by comparing model-based and census-based growth rates (Coulson *et al.*, 2001) or by checking visually the similarity of model-based and census-based trajectories (Kanyamibwa and Lebreton, 1992).

Estimates of the annual survival probabilities of wild animals are usually obtained from appropriate mark–recapture–recovery experiments. In these experiments, animals receive marks which identify them uniquely. Subsequent observations of these marked animals, either alive or dead, result in data which are fitted by probability models incorporating annual probabilities of survival. The models that are used can be quite sophisticated, involving age and time

Address for correspondence: B. J. T. Morgan, Institute of Mathematics and Statistics, Cornwallis Building, University of Kent at Canterbury, Canterbury, Kent, CT2 7NF, UK.
E-mail: B.J.T.Morgan@ukc.ac.uk

dependence of parameters as well as dependence on individual and environmental covariates. An illustration is provided by Catchpole *et al.* (2000).

Besbeas *et al.* (2002) adapted Kalman filter theory to form a single combined likelihood, with different components corresponding to survival and census data, and sharing common parameters. Maximization of the combined likelihood then produces parameter estimates which describe all the data simultaneously. This approach centres on a state space model for the census data with a likelihood derived by a Kalman filter and provides a natural synthesis of models for mark–recapture–recovery data and models for population processes.

However, a drawback of this approach is that it requires specialized computer code, not only for the Kalman filter component, but also for the additional likelihood components, some of which may be complex, and may even have been derived by using specialist computer packages. This problem is likely to preclude the use of the approach of Besbeas *et al.* (2002) in practice, although we believe that it is likely to become an important tool for many studies of animal populations. In this paper we suggest a way of overcoming this difficulty, by means of suitable multivariate normal approximations. We evaluate the effectiveness of the approach by applying it to the two examples examined in detail in Besbeas *et al.* (2002).

2. Data, models and approximation

We shall illustrate our approach by using two examples, in each of which ring–recovery data from marked birds are combined with abundance data. However, the approach is quite general, as we shall see. The two examples that we consider are described in Besbeas *et al.* (2002) and are related to observations on lapwings and herons. As explained in Besbeas *et al.* (2002) these are two extensive and important data sets. For common lapwings, *Vanellus vanellus*, breeding in Britain, the Common Birds Census index (Taylor, 1965; Fewster *et al.*, 2000) provides an indication of the national breeding population, from 1965 to 1998 inclusively, whereas the ring–recovery data provide the numbers of birds recovered dead in successive years after being ringed as chicks from 1963 to 1997. For British grey herons, *Ardea cinerea*, census data estimate the numbers of breeding pairs in England and Wales between 1928 and 1998 inclusively, and the ring–recovery data are for birds ringed in the nest between 1955 and 1997. The raw data are given in Besbeas *et al.* (2001). As an illustration, we show in Table 1 the lapwing recovery data from 1990 to 1997.

Discussion of model formation and selection from a classical perspective is provided in Besbeas *et al.* (2002), using Akaike's information criterion and incorporating ecological knowledge. The Bayesian perspective is given by Brooks *et al.* (2000, 2002). We model the recovery data by using annual survival probabilities ϕ , with components which describe age dependence, and a recovery probability λ , which is the probability of recovery and reporting of marked dead birds, and which varies over time. For the lapwings there are two age classes of survival, corresponding to birds in their first year of life and older birds, whereas for the herons there are three age classes of survival, as we also distinguish a separate probability of survival for birds in their second year of life in that case. All the annual survival probabilities are regressed on a single measure of winter severity, w , using logistic regression. The covariate w gives the number of days in the year that a measure of central England temperature fell below freezing. Thus we have, for both species, $\text{logit}(\phi_1) = \beta_0 + \beta w$, for herons only, $\text{logit}(\phi_2) = \gamma_0 + \gamma w$, and $\text{logit}(\phi_a) = \delta_0 + \delta w$, which applies to birds aged 1 year or older for lapwings and 2 years or older for herons. In addition, the reporting rates for dead wild birds are generally found to be declining over time, and so we set $\text{logit}(\lambda) = \nu_0 + \nu t$, where t measures years. We do not here consider overdispersion, but that is easily incorporated by means of suitable additive random effects, as described by Barry *et al.* (2002).

Table 1. Illustrative recovery subtable for lapwings, taken from Besbeas *et al.* (2001)[†]

Year of ringing	Number ringed	Numbers in the following years of recovery:								Number of days below freezing (<i>w</i>)
		1991	1992	1993	1994	1995	1996	1997	1998	
1990	4170	12	3	3	2	1	0	2	0	12
1991	4314		9	4	6	1	0	1	0	12
1992	3480			18	3	1	2	0	1	9
1993	3689				6	5	2	2	1	6
1994	3922					12	4	6	0	3
1995	3591						7	5	1	18
1996	4488							7	0	10
1997	4339								5	0

[†]The data show the numbers of British lapwings recovered dead in successive years after being ringed as chicks, for birds ringed between 1990 and 1997 inclusively. Reporting rates of dead birds have been declining over time. At the start of the lapwing study similar numbers of birds found dead are obtained from about 50% of the cohort sizes of ringed birds shown. Also shown is the number *w* of days in the year that a measure of central England temperature fell below freezing. Note that by a year *i*, say, we mean the period of time from April 1st of year *i* to March 31st of year *i* + 1.

The census data are described by means of a state space model based on a Leslie matrix, and involving a productivity measure *p* and measurement error variance σ^2 , in addition to the survival probabilities. For the lapwings, of which numbers are declining, we set $\log(p) = \kappa_0 + \kappa t$, but for herons *p* is taken as constant, $\log(p) = \kappa_0$. Maximum likelihood parameter estimates result in general from maximizing the joint likelihood

$$L_j(\phi, \lambda, p, \sigma) = L_r(\phi, \lambda) L_c(\phi, p, \sigma). \quad (1)$$

The census component $L_c(\phi, p, \sigma)$ is formed by a Kalman filter, whereas the ring-recovery component $L_r(\phi, \lambda)$ is a product of multinomial distributions, with each corresponding to a separate cohort of ringed birds. For the lapwing example, if $N_{1,t}$ denotes the number of female birds that are 1 year old at time *t*, and $N_{a,t}$ denotes the number of female birds aged 2 years or older at time *t*, the state space transition equation used by Besbeas *et al.* (2002) is

$$\begin{pmatrix} N_1 \\ N_a \end{pmatrix}_t = \begin{pmatrix} 0 & p\phi_1 \\ \phi_a & \phi_a \end{pmatrix} \begin{pmatrix} N_1 \\ N_a \end{pmatrix}_{t-1} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_a \end{pmatrix}_t.$$

This assumes that breeding starts at age 2 years. We assume also that only breeding birds are censused, so that the annual indices $\{y_t\}$ are related to the population level by the measurement equation

$$y_t = (0 \quad 1) \begin{pmatrix} N_1 \\ N_a \end{pmatrix}_t + \eta_t.$$

Assumptions regarding the distribution and variances of $\{\varepsilon_{1,t}\}$, $\{\varepsilon_{a,t}\}$ and $\{\eta_t\}$ are explained in Besbeas *et al.* (2002). A primary advantage of combining the data sets is that it allows us to estimate *p*, which is not possible from the census data alone.

The independence assumption that is made in equation (1) is unlikely to be seriously violated in practice. In this paper, we make a multivariate normal approximation to $L_r(\phi, \lambda)$, i.e.

$$2 \log\{L_r(\phi, \lambda)\} = \text{constant} - (\hat{\theta} - \theta)' \hat{\Sigma}^{-1} (\hat{\theta} - \theta),$$

where we write θ to denote the model parameters on the logistic scale, and where $\hat{\theta}$ and $\hat{\Sigma}$ are respectively the maximum likelihood estimates of θ and the dispersion matrix of $\hat{\theta}$, obtained from a separate model fitting exercise for the ring–recovery data alone. This approach is motivated by the asymptotically multivariate normal distribution of maximum likelihood estimators and has been used to good effect in Lebreton *et al.* (1995). It clearly greatly simplifies both the resulting form for L_j and its maximization.

3. Results

For the two examples of the last section we evaluate the approximation in Table 2; the standard errors are obtained from the observed information matrices. The agreement between the exact and approximate results is seen to be very good. An additional benefit from using the approximate approach is that, in comparison with the exact analysis, it is far less sensitive to starting values for the maximum likelihood iteration. This is especially important if we are forming bootstrap estimates of error. For example, for the simpler lapwing example, from 100 random starting values taken near the solution, the divergence rate of the approximate method was 10%, compared with 21% for the exact method, when a quasi-Newton iterative procedure was used to maximize the likelihood. Additionally, we observed that for the examples of this paper the approximate analysis was about 2.5 times faster than the exact method. Absolute computation times for the exact joint analysis on a Sun Enterprise 400 MHz 250 computer were 233 s for

Table 2. Evaluation of the multivariate normal approximation for the ring–recovery likelihood, for lapwings and herons†

Parameter	Ring–recovery alone	Exact combination	Approximate combination
<i>Lapwings</i>			
β_0	0.5158 (0.0679)	0.5231 (0.0679)	0.5226 (0.0678)
β	−0.0241 (0.0072)	−0.0228 (0.0070)	−0.0227 (0.0070)
δ_0	1.5011 (0.0685)	1.5210 (0.0693)	1.5191 (0.0686)
δ	−0.0360 (0.0051)	−0.0279 (0.0045)	−0.0280 (0.0045)
ν_0	−4.5668 (0.0351)	−4.5632 (0.0352)	−4.5634 (0.0351)
ν	−0.5729 (0.0640)	−0.5841 (0.0637)	−0.5837 (0.0638)
κ_0		−1.1513 (0.0886)	−1.1489 (0.0876)
κ		−0.4323 (0.0743)	−0.4314 (0.0740)
σ		159.4691 (22.0625)	159.6127 (21.8749)
<i>Hérons</i>			
β_0	−0.2024 (0.0480)	−0.1908 (0.0478)	−0.1901 (0.0479)
β	−0.0309 (0.0054)	−0.0217 (0.0048)	−0.0215 (0.0048)
γ_0	0.3745 (0.0730)	0.3842 (0.0734)	0.3837 (0.0728)
γ	−0.0155 (0.0068)	−0.0166 (0.0060)	−0.0166 (0.0060)
δ_0	1.1655 (0.0706)	1.1871 (0.0708)	1.1861 (0.0700)
δ	−0.0110 (0.0052)	−0.0140 (0.0038)	−0.0140 (0.0038)
ν_0	−2.0309 (0.0254)	−2.0294 (0.0254)	−2.0295 (0.0254)
ν	−0.8302 (0.0462)	−0.8337 (0.0460)	−0.8335 (0.0461)
κ_0		−0.0459 (0.0715)	−0.0455 (0.0708)
σ		464.7502 (43.4774)	464.2624 (43.7103)

†In each case we show the maximum likelihood parameter estimate and the corresponding standard error. The estimates of error are obtained by a second-order finite difference approximation applied to the log-likelihood.

herons and 112 s for lapwings. Thus here the time savings are small, but for more complex examples, such as that mentioned in Section 4, they are appreciable.

We can see from Table 2 that the shift in $\hat{\theta}$ from the recovery analysis alone to the joint analysis is not large. It is both the overall magnitude of this change and the size of the data sets which determine the effectiveness of the multivariate normal approximation that is made here. This is because we require the multivariate normal approximation to L_r to be good for the value of θ which maximizes L_j , and not just for the value that maximizes L_r . We expect that in practice for most examples the approximation will behave well. We show in Fig. 1 the good agreement between the observed and expected $\log\{L_r(\phi, \lambda)\}$ for the small illustrative data set of Table 1, by means of profile log-likelihoods, in this case for the parameter β_0 . The value $\hat{\beta}_0 = 0.1786$ differs from $\hat{\beta}_0 = 0.5231$, given in Table 2, since the data of Table 1 are atypical when compared with the complete, 35-year ring-recovery data set. This good agreement results for profiles with respect to other parameters, and also for both the complete data sets. Without conducting an extensive simulation study, it is not possible to be clear when the approximation of the paper might fail in practice. To obtain further evidence for the good performance of the approximation we have sampled ring-recovery data sets from the full ring-recovery set for both herons and lapwings. Thus, for the lapwing data set, 26 such (10×10) -year ring-recovery data sets were obtained, for years 1965–1974, 1966–1975, etc., and, for the longer heron data set, 29 (15×15) -year ring-recovery data sets were similarly obtained. For each small data set we carried out a ring-recovery analysis alone, an exact joint analysis and an approximate joint analysis, with both joint analyses making use of all the census data. Comparisons were made by using the Mahalanobis distance, and these are presented in Table 3. With only one exception for each species, the discrepancy between the two joint an-

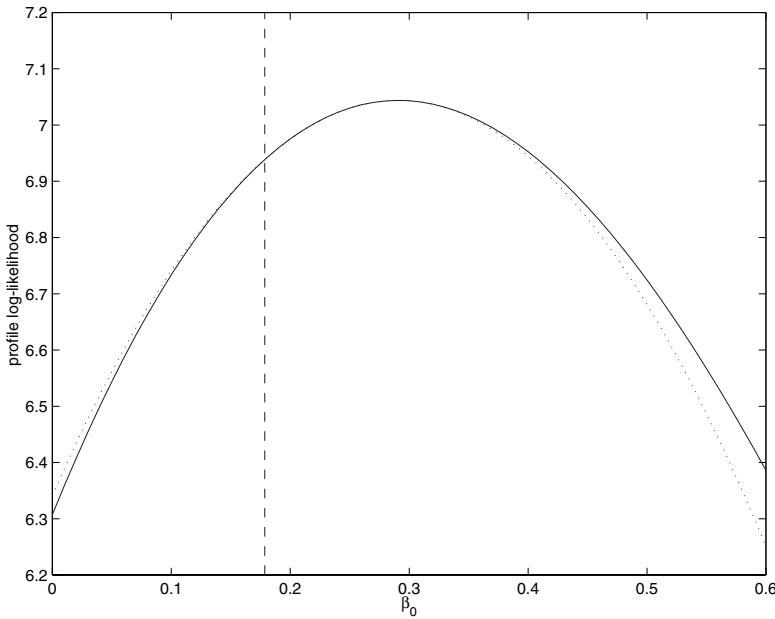


Fig. 1. Agreement between $\log\{L_r(\phi, \lambda)\}$ and the multivariate normal approximation for the data of Table 1: the two curves are profile log-likelihoods taken with respect to the parameter β_0 (— exact; approximate; |, location of $\hat{\beta}_0 = 0.1786$, the value which maximizes the combined exact likelihood, making use of the entire run of the census data for lapwings)

Table 3. Mahalanobis distances between sets of parameter estimates, for lapwing and heron subsets of ring–recovery data†

Data set	Distances for lapwings		Distances for herons		Data set	Distances for lapwings		Distances for herons	
	(a)	(b)	(a)	(b)		(a)	(b)	(a)	(b)
1	0.03	1.06	0.68	4.48	16	0.29	1.44	0.11	2.61
2	0.07	1.35	0.55	4.06	17	0.64	1.81	0.13	2.57
3	0.11	0.82	0.52	3.88	18	0.45	1.34	0.05	2.05
4	0.08	1.69	0.51	3.74	19	0.06	2.26	0.05	2.09
5	0.14	1.48	0.49	3.69	20	0.12	2.15	0.08	2.32
6	0.10	1.29	3.94	3.80	21	0.08	2.19	0.09	2.31
7	1.80	1.30	0.51	3.09	22	0.04	1.25	0.13	2.64
8	0.21	0.73	0.35	3.30	23	0.02	1.00	0.12	2.54
9	0.08	0.47	0.09	2.91	24	0.04	1.33	0.11	2.53
10	0.18	0.96	0.66	4.21	25	0.11	1.63	0.10	2.42
11	0.08	1.04	0.41	3.78	26	0.09	1.23	0.10	2.29
12	0.05	0.75	0.19	3.10	27			0.09	2.37
13	0.07	0.44	0.05	1.76	28			0.09	2.18
14	0.03	0.48	0.06	2.08	29			0.09	2.40
15	0.12	0.84	0.04	1.80					

†(a) Distance between estimates from the exact joint analysis and approximate joint analysis, using the estimate of the variance–covariance matrix obtained from the exact joint analysis; (b) distance between estimates from the ring–recovery data set alone and approximate joint analysis, using the estimate of the variance–covariance matrix obtained from the ring–recovery data set. Three fewer parameters are involved in (b) compared with (a) for lapwings, and two fewer for herons.

alyses is seen to be appreciably smaller than the discrepancy between the ring–recovery alone set of parameter estimates and the approximate joint set of parameter estimates.

In practice, a simple check of the adequacy of the approximation lies in the goodness of fit of the model, following the combined analysis. A combined analysis should follow a good fit to the ring–recovery data alone. If the fit of the combined model is poor in some respects, then either the approximation is at fault or the assumption of common parameters for the various components of the analyses is incorrect.

4. Discussion

The very good performance of the approximation that is reported here is encouraging. It will greatly facilitate the use of the combined analysis procedure of Besbeas *et al.* (2002). When we assemble combined likelihoods then we can envisage doing this for additional different aspects of the model—we might for example introduce a component measuring fecundity, as suggested for instance by Pollock and Cornelius (1988), if appropriate data are available. The approximate approach would remove the need to produce, in that case, a separate specialized maximum likelihood estimation program for data on nesting success. All that is needed are the appropriate fecundity parameter estimates $\hat{\theta}$ and $\hat{\Sigma}$ from a previous analysis. There are implications here for the reporting of results from analyses that might later be used in multivariate normal approximations, since typically what are published are only $\hat{\theta}$ and $\text{diag}(\hat{\Sigma})$. A particular attraction of the multivariate normal approximation is that if, as in the examples of this paper, the primary

interest lies in ϕ , then we need only $\hat{\phi}$ and an estimate of its dispersion matrix, as shown in Appendix A. This result will naturally produce a further increase in computational efficiency. The only possible defect of the approach of this paper is that it would not allow for experimentation with more elaborate models, if that should become feasible as a result of combining separate data sets. We have used the approach of this paper in complex modelling examples, as in Catchpole *et al.* (2000), for which census and fecundity data are also available. The model synthesis that was achieved was only feasible as a result of using the approximation of the paper.

Acknowledgements

We thank the very many volunteer bird ringers and census workers whose dedication has made the study of this paper possible, and Stephen Freeman and Ted Catchpole for their input to the state space modelling which underpins the work of this paper. The work of PB was supported by Biotechnology and Biological Science Research Council–Engineering and Physical Sciences Research Council grant 96/E09745.

Appendix A

Suppose that we write

$$L(\mu) = (2\pi)^{-k/2} |\Sigma|^{-1/2} \exp\{-\frac{1}{2}(\theta - \mu)' \Sigma^{-1}(\theta - \mu)\},$$

where $k = \dim(\theta)$, and we similarly partition $\theta' = (\theta_1, \theta_2)$ and $\mu' = (\mu_1, \mu_2)$, regarding μ_1 as the parameters of interest. Let $k_1 = \dim(\theta_1) \equiv \dim(\mu_1)$. We suppose also that the corresponding partition of Σ^{-1} is written as

$$\Sigma^{-1} = \begin{pmatrix} \Sigma^{11} & \Sigma^{12} \\ \Sigma^{21} & \Sigma^{22} \end{pmatrix}.$$

We are given θ and Σ , and we require the maximum likelihood estimate of μ_1 . We can write $l(\mu) = \log\{L(\mu)\}$ as

$$l(\mu) = \text{constant} - \frac{1}{2}(\theta_1 - \mu_1)' \Sigma^{11}(\theta_1 - \mu_1) - \frac{1}{2}(\theta_1 - \mu_1)' \Sigma^{12}(\theta_2 - \mu_2) \\ - \frac{1}{2}(\theta_2 - \mu_2)' \Sigma^{21}(\theta_1 - \mu_1) - \frac{1}{2}(\theta_2 - \mu_2)' \Sigma^{22}(\theta_2 - \mu_2).$$

From setting $\partial l / \partial \mu_1 = \partial l / \partial \mu_2 = 0$, we obtain the two equations

$$\begin{aligned} \Sigma^{11}(\theta_1 - \mu_1) + \Sigma^{12}(\theta_2 - \mu_2) &= 0, \\ \Sigma^{21}(\theta_1 - \mu_1) + \Sigma^{22}(\theta_2 - \mu_2) &= 0. \end{aligned} \quad (2)$$

Eliminating $\theta_2 - \mu_2$ then gives

$$\{\Sigma^{11} - \Sigma^{12}(\Sigma^{22})^{-1}\Sigma^{21}\}(\theta_1 - \mu_1) = 0,$$

i.e.

$$\Sigma_{11}^{-1}(\theta_1 - \mu_1) = 0, \quad (3)$$

where Σ_{11} is the $k_1 \times k_1$ top left-hand submatrix of Σ . Normally, when likelihoods are combined, the right-hand sides of equations (2) and (3) will be functions of μ_1 . We therefore see that the approximation of this paper only requires estimates of θ_1 and the corresponding dispersion matrix Σ_{11} .

References

- Barry, S. C., Brooks, S. P., Catchpole, E. A. and Morgan, B. J. T. (2002) The analysis of ring-recovery data using random effects. *Biometrics*, to be published.
- Besbeas, P., Freeman, S. N., Morgan, B. J. T. and Catchpole, E. A. (2001) Stochastic models for animal abundance and demographic data. *Technical Report UKC/IMS/01/16*. University of Kent at Canterbury, Canterbury.

- Besbeas, P., Freeman, S. N., Morgan, B. J. T. and Catchpole, E. A. (2002) Integrating mark-recapture-recovery and census data to estimate animal abundance and demographic parameters. *Biometrics*, **58**, 540–547.
- Brooks, S. P., Catchpole, E. A. and Morgan, B. J. T. (2000) Bayesian animal survival estimation. *Statist. Sci.*, **15**, 357–376.
- Brooks, S. P., Catchpole, E. A., Morgan, B. J. T. and Harris, M. P. (2002) Bayesian methods for analysing ringing data. *J. Appl. Statist.*, **29**, 187–206.
- Caswell, H. (2000) *Matrix Population Models*. Sunderland: Sinauer.
- Catchpole, E. A., Morgan, B. J. T., Coulson, T. N., Freeman, S. N. and Albon, S. D. (2000) Factors influencing Soay sheep survival. *Appl. Statist.*, **49**, 453–472.
- Coulson, T. N., Catchpole, E. A., Albon, S. D., Morgan, B. J. T., Pemberton, J. M., Clutton-Brock, T. H., Crawley, M. J. and Grenfell, B. T. (2001) Age, sex, density, winter weather and population crashes in Soay sheep. *Science*, **292**, 1528–1531.
- Fewster, R. M., Buckland, S. T., Siriwardena, G. M., Baillie, S. R. and Wilson, J. D. (2000) Analysis of population trends for farmland birds using generalized additive models. *Ecology*, **81**, 1970–1984.
- Kanyamibwa, S. and Lebreton, J.-D. (1992) Variation des effectifs de Cigogne blanche et facteurs de milieu: un modèle démographique. In *Les Cigognes d'Europe* (eds J.-L. Mériaux, A. Schierer, C. Tombal and J.-C. Tombal), pp. 259–264. Metz: Institut Européen d'Ecologie.
- Lebreton, J.-D., Morgan, B. J. T., Pradel, R. and Freeman, S. N. (1995) A simultaneous survival rate analysis of dead recovery and live recapture data. *Biometrics*, **51**, 1418–1428.
- Leslie, P. H. (1945) On the use of matrices in certain population mathematics. *Biometrika*, **33**, 183–212.
- Pollock, K. H. and Cornelius, W. L. (1988) A distribution-free nest survival model. *Biometrics*, **44**, 397–404.
- Taylor, S. M. (1965) The Common Birds Census—some statistical aspects. *Bird Stud.*, **12**, 268–286.