# Communications in Statistics - Simulation and Computation

## Selection of Models of Lagged Identification Rates and Lagged Association Rates Using AIC and QAIC

Hal Whitehead [a]

[a] Department of Biology, Dalhousie University, Halifax, Nova Scotia, Canada

PLEASE SCROLL DOWN FOR ARTICLE

Taylor & Francis
Taylor & Francis Group

# Inference

# Selection of Models of Lagged Identification Rates and Lagged Association Rates Using AIC and QAIC

## HAL WHITEHEAD

Department of Biology, Dalhousie University, Halifax,
Nova Scotia, Canada

*The lagged identification rate is the probability of identifying an individual given its identification some time lag earlier. The lagged association rate is the probability that two individuals are associated given their association some time lag earlier. Models of lagged identification and association rates fit by maximizing the sums of non independent log-likelihoods have approximately unbiased parameter estimates. Simulations suggest that: Akaike-Information-Criterion often selects the true model of lagged identification rate data; quasi-Akaike-Information-Criterion performs better for lagged association rates; and confidence intervals for parameters are best obtained by bootstrap methods for lagged identification rates and quasi-likelihood or jackknife methods for lagged association rates.*

## 1. Introduction

In recent years, selection among models being fitted to biological data has progressed to become an important subdiscipline of biometrics (Burnham and Anderson, 2002). The Akaike Information Criterion (AIC; Akaike, 1973), and variants of it, have become widely used in selecting appropriate models. For instance, computational routines for mark-recapture methods of estimating population parameters from data on identifications of marked individuals now routinely include the AIC for model selection or model averaging (e.g., Pledger et al., 2003; White and Burnham, 1999).

Records of identifications of individuals can give insight into other areas of biology, including movements (e.g., Hilborn, 1990) and social structure (e.g., Bejder et al., 1998). We seek mathematical models of these phenomena from identification

data, analogous to those that have been achieved for population analyses. Records of individual identifications often can be used for movement and social analysis, as well as population assessment. Therefore, the raw data consist of records of who was where, when, and with whom.

In the analysis of movement among subareas of a habitat, a very useful output is a model estimating the probability that an individual in area A is the same as an individual identified in area A (or some other area B) after some time lag, $\tau$, a function which the author has called the lagged identification rate ($R(\tau)$; Whitehead, 2001). Similarly, in social analyses we can seek models of lagged association rates, the probability that if two individuals are associated now, they will still be associated $\tau$ time units later ($g(\tau)$; Whitehead, 1995). A variant of this, the standardized lagged association rate, $g'(\tau)$, is the probability that if $X$ and $Y$ are associates, then $\tau$ time units later a randomly-chosen associate of $X$ is $Y$. Standardized lagged association rates are appropriate measures in cases when all associates may not be identified (Whitehead, 1995).

Log-likelihoods are the basis of most methods of model selection (Burnham and Anderson, 2002), but unfortunately, their direct calculation is challenging with even quite small data sets and fairly simple models of movement or sociality, and becomes practically impossible for data that have large time spans and/or more complex models. For instance, to calculate the likelihood that, during 30 sampling periods, an individual is identified in area A during periods 1, 12, and 20 under a simple model of migration between a study area and some other habitat involves considering the probability of $2^{27} = 134{,}217{,}728$ possible movement histories (as the individual could have been in either of the areas in each of the 27 periods during which it was not identified). With social models in which the behavior of each member of a dyad (pair of individuals) must be considered, the difficulties escalate even more rapidly.

One solution to this challenge is to consider the sum of log-likelihood elements which are not necessarily independent (Whitehead, 1995, 2001). In practice, the author recommends calculating the log-likelihood of the observed reidentifications between each pair of sampling periods, ignoring the data for all other periods, and then summing these. So, for lagged identification rates, the calculated quantity is:

$$L^* = \sum_\tau \sum_{j,k \,|\, (t_j - t_k) = \tau} \mathbf{Log}\left(L(R(\tau) \,|\, m_{jk}, n_j, n_k)\right), \tag{1}$$

where $t_j$ is the time of sampling period $j$, $n_j$ is the number of individuals identified in period $j$, $m_{jk}$ is the number of individuals identified in both periods $j$ and $k$, and $L(R(\tau) \,|\, m_{jk}, n_j, n_k)$ is the likelihood of the model $R(\tau)$ given the identification and reidentification data for each pair of sampling periods.

For lagged association rates, the formulation is:

$$L^* = \sum_\tau \sum_{j,k \,|\, (t_j - t_k) = \tau} \mathbf{Log}(L(g(\tau) \,|\, \{a_j(X, Y)\}, \{a_k(X, Y)\})), \tag{2}$$

where $a_j(X, Y) = 1$ if $X$ and $Y$ were recorded as associated in time period $j$, and $a_j(X, Y) = 0$ if they were not associated, or either was not identified, during the sampling period.

In Eqs. (1) and (2), the pairs of periods are grouped by inter-period lag, $\tau$, as this is the dependent variable of the models, and this formulation simplifies calculation.

These two-period likelihoods are each calculated easily, even under quite complex models (and their calculation can be simplified in various ways; see Whitehead, 2001), but they are not independent. The author has shown analytically and using simulation that maximizing the sum of the non independent log-likelihoods, $L^*$, produces approximately unbiased estimates of model parameters (Whitehead, 2001, Appendix), but could not justify their use in estimating precision or in model selection. Despite this, likelihood-based models have been used to choose between models of lagged association and identification rates (e.g., Gowans et al., 2000; Karczmarski et al., 2005; Ottensmeyer and Whitehead, 2003). The author is partially responsible for this as he did not initially appreciate the non independence problem (Whitehead, 1995), and included likelihood, and later AIC, in his computer program SOCPROG (http://myweb.dal.ca/~hwhitehe/social.htm) which calculates lagged identification and association rates.

Bootstrap and jackknife techniques can be used to estimate parameter precision for lagged identification rates (Whitehead, 2001). For lagged association rates, bootstrap methods are invalid as resampling the same individual will indicate more social stability than is real. Jackknife techniques are usable (Whitehead, 1995), but tend to be both conservative and approximate (Efron and Stein, 1981).

Here we use simulation to explore two possibilities: (a) that despite the non independence of the summed log-likelihoods, AIC might still give useful guidance in selecting models of lagged identification rates and lagged association rates; (b) that the quasi-likelihood variant of AIC, QAIC, which compensates for overdispersed count data when using Poisson or binomial models (Burnham and Anderson, 2002), as in the lagged identification rate and lagged association rate models, might compensate for the summing of non independent log-likelihoods. Secondly, we compare bootstrap, jackknife, likelihood, and quasi-likelihood methods of calculating confidence intervals of estimated parameters.

## 2. Methods

### 2.1. *Models of Lagged Identification Rates*

Three realistic types of population were simulated, using MATLAB 6.5 (Mathworks, Natick, MA), over 605 consecutive arbitrary time units:

A: A closed population of $N = 100$ individuals present in the study area throughout with no birth, death, immigration or emigration;

B: A population of $N = 100$ individuals in the study area with permanent emigration at a rate of $\lambda = 0.008$ per individual per time unit, with departed individuals being replaced 1:1 by new individuals;

C: A closed population of $Z = 300$ individuals, members of which can be either inside or outside the study area. Individuals in the study area leave the study area at a rate of $\lambda = 0.08$ per individual per time unit and individuals outside the study area reenter it with a probability of $\mu = 0.04$ per individual per time unit. In this case there will be approximately $N = 100$ individuals in the study area at any time (with expected numbers entering and leaving equal at 8 individuals per time unit).

Each population was randomly sampled at 25 time periods, $t_j = 1$–5, 51–55, 101–105, 501–505, 601–605, indicating irregularly spaced field effort, a common

scenario. One hundred populations were simulated for each of the three population types and each of 10, 20, or 40 randomly-chosen individuals from the study area identified per sampling period, as well as with $N = 1,000$, or (for model C) $Z = 3,000$. The lagged identification rates were estimated from the identification data produced from each simulated population using (Whitehead, 2001):

$$R(\tau) = \frac{\sum_{j,k \,|\, (t_k - t_j) = \tau} m_{jk}}{\sum_{j,k \,|\, (t_k - t_j) = \tau} n_j \cdot n_k}. \tag{3}$$

Three models were fitted to these data using likelihood methods assuming that each $m_{jk}$ was binomially distributed with parameters $n_j \cdot n_k$ and $R(t_k - t_j)$, as in Whitehead (2001), even though the data are neither independent nor binomially distributed. The models were (as in Gowans et al., 2000):

$$1: \ R(\tau) = \alpha$$
$$2: \ R(\tau) = \alpha \cdot e^{-\beta \cdot \tau}$$
$$3: \ R(\tau) = \gamma \cdot e^{-\beta \cdot \tau} + \alpha.$$

Population Type A theoretically fits model 1 as there is no autocorrelation; population Type B, a Poisson process, theoretically fits model 2 with $\alpha = 1/N$ and $\beta = \lambda$; population Type C theoretically fits model 3 with $\beta = (\lambda + \mu)$, $\alpha = \mu/((\lambda + \mu) \cdot N)$, $\gamma = \lambda/((\lambda + \mu) \cdot N)$ (for derivations of these expressions see Whitehead, 2001).

To set the results in perspective, the author constructed a series of data sets as in populations of Type B but with true independent binomial structure by replacing the $m_{jk}$'s with random numbers from binomial distributions with parameters $n_j \cdot n_k$ and $\alpha \cdot e^{-\beta \cdot (t_k - t_j)}$ (model 2) and $\alpha = 0.01$ and $\beta = 0.008$ as in theory for population Type B. We will call these data sets, which have a true independent binomial structure, population Type BX.

## 2.2. Models of Lagged Association Rates

Three realistic types of social structure were simulated, all in a closed population of $N = 100$ individuals:

D: Individuals form $W$ groups of mean size 10 (so $W = N/10 = 10$), with random allocation of individuals to groups at each sampling period-random associations;

E: Individuals form $W$ groups of mean size 10 (so $W = N/10 = 10$), and change groups with probability $\lambda = 0.008$ per time unit-casual acquaintances;

F: Individuals are randomly allocated to $U = 20$ permanent social units, and these are randomly allocated to $W = 10$ groups, with units changing groups with probability $\lambda = 0.008$ per time unit-permanent companions plus casual acquaintances.

The sampling schemes were as with the lagged identification rates (above), except that 1, 2, or 4 groups were sampled each period (rather than 10, 20, or 40 individuals) and all members of each selected group were identified. Lagged

association rates, $g(\tau)$, were estimated from the identification data produced from each simulated population using (Whitehead, 1995):

$$g(\tau) = \frac{\sum_{j,k \mid (t_k - t_j) = \tau} \sum_X \sum_{Y \neq X} a_j(X, Y) \times a_k(X, Y)}{\sum_{j,k \mid (t_k - t_j) = \tau} \sum_X \sum_{Y \neq X} a_j(X, Y)}. \tag{4}$$

In this expression both the numerator and denominator are summed over all pairs of sampling periods $\tau$ time units apart. Summing over each such pair, and over all individuals $(X)$, the numerator is the number of other individuals $(Y)$ identified in the same group as $X$ in both period $j$ and period $k$, while the denominator is the number of other individuals identified in the same group as $X$ in period $j$. $g(\tau)$ is an estimate of the probability that if $X$ and $Y$ are identified in the same group, then, $\tau$ time units later, $Y$ is once again in the same group as $X$. If group membership is random, $g(\tau)$ is approximately $(1/W)$; if individuals form permanent groups, then $g(\tau) = 1$.

The following three models of lagged association rates (from Whitehead, 1995) were fit to the simulated data using a binomial model as with the lagged identification rates:

$$1: \ g(\tau) = \alpha$$
$$4: \ g(\tau) = (1 - \alpha) \cdot \mathrm{e}^{-\beta \cdot \tau} + \alpha$$
$$5: \ g(\tau) = ((1 - \alpha) \cdot \mathrm{e}^{-\beta \cdot \tau} + \alpha) \cdot \mathrm{e}^{-\gamma \cdot \tau}$$

Models 2 and 3 are not very realistic for lagged association rate data, as there would be a finite possibility of individuals disassociating over infinitely short times ($\tau = 0$). Models 4 and 5 address this structurally by making $g(\tau) = 1$ when $\tau = 0$, and model 1 omits all mention of time lag.

Population Type D, theoretically fits model 1 with $\alpha = 1/W$; population Type E theoretically fits model 4 with $\alpha = 1/W$ and $\beta = 2 \cdot \lambda \cdot W/(W - 1)$; population Type F also theoretically fits model 4 with $\alpha = (W^2 + U - 1)/[W \cdot (W + U - 1)]$ and $\beta = 2 \cdot \lambda \cdot W/(W - 1)$ (derivations in Appendix). Model 5 introduces mortality or permanent emigration, potentially realistic features, but not present in the simulated data.

As with the lagged identification rates, the author constructed a series of data sets with true binomial structure by replacing the numerator of Eq. (4) in populations of Type E with random numbers from binomial distributions with parameters $\sum \sum a_j(X, Y)$ and $(1 - \alpha) \cdot \mathrm{e}^{-\beta \cdot (t_k - t_j)} + \alpha$ (model 4) where $\alpha = 1/W = 0.1$ and $\beta = 2 \cdot \lambda \cdot W/(W - 1) = 0.0178$ as in theory for population Type E. We will call these data sets, which have a true independent binomial structure, population Type EX.

### 2.3. *Output from Model Fitting*

Output from fitting each model to each random data set includes the estimates of parameters $(\alpha, \beta, \gamma)$ from maximizing the (summed) log-likelihoods, the maximum (summed) log-likelihoods ($L^*$, from Eqs. (1) and (2)), and an estimate of the variance inflation factor, $c$. $c$ is estimated from the ratio of the goodness-of-fit $\chi^2$-statistic to its degrees of freedom, $v$ (Burnham and Anderson, 2002). The $\chi^2$-statistics were calculated by comparing the observed total number of pairs of identifications $\tau$ time

units apart which were of the same individual with the expected number given the model and its estimated parameters. Categories were lumped so that all time lag categories contained an expected value of at least six. The degrees of freedom, $v$, was calculated as the number of time delay categories minus the number of parameters in the model, $K$, minus one. Two model fitting criteria were considered (Burnham and Anderson, 2002):

$$\text{AIC} = -2 \cdot L^* + 2 \cdot K$$
$$\text{QAIC} = -2 \cdot L^*/\hat{c} + 2 \cdot K.$$

For each simulated data set, $\hat{c}$ was $c$ as calculated for the most general model. If $\hat{c} > 1$, the QAIC is the preferable criterion (Burnham and Anderson, 2002); if $\hat{c} < 1$ then the QAIC reverts to the AIC (i.e., $\hat{c} = 1$). The model with the smallest AIC or QAIC was selected as the best-fitting for that data set and criterion.

There are also second-order versions of AIC and QAIC, which correct for small sample sizes (Burnham and Anderson, 2002), $\text{AIC}_c$ and $\text{QAIC}_c$. The author examined their performance on the data sets with smallest sample sizes (ten individuals per sample with lagged identification rates, and one group per sample with lagged association rates), using the formulae:

$$\text{AIC}_c = -2 \cdot L^* + 2 \cdot K \cdot M/(M - K - 1)$$
$$\text{QAIC}_c = -2 \cdot L^*/\hat{c} + 2 \cdot K \cdot M/(M - K - 1),$$

where $M$ is the sample size (sum, over all lags, of the numerators of Eq. (3) for lagged identification rates, or Eq. (4) for lagged association rates).

### 2.4. *Confidence Intervals of Parameter Estimates*

For each population type, the author compared the performance of several techniques of estimating the precision of parameter estimates, using only 20 identifications per sampling interval and the theoretically correct model of lagged identification or association rates for each population type. The true span of the 95% confidence interval was estimated from the distribution of the parameter estimates of the 100 runs of each population type, and for 100 additional runs 95% confidence intervals were estimated for each parameter in four ways:

- bootstrap using 1,000 samples, with replacement, of the identified individuals (only for lagged identification rates);
- jackknife with sampling periods within 50-time-unit intervals being omitted in turn (Whitehead, 1995), giving five jackknife replicates;
- likelihood support interval, defined for any parameter as the range of values such that the (summed) log-likelihood, maximizing over all other parameters, is within two of the overall maximum (summed) log-likelihood (Edwards, 1992);
- quasi-likelihood support interval, defined for any parameter as the range of values such that the (summed) log-likelihood, maximizing over all other parameters, is within $2 \cdot \hat{c}$ of the overall maximum (summed) log-likelihood.

For each population type and parameter of the correct model, we present the median and standard deviation of the spans of the estimated 95% confidence intervals

over the model runs, as well as the coverage, the percentage of confidence intervals containing the true parameter value. The coverage should theoretically be 95%.

## 3. Results

### 3.1. *Model-Fitting: Lagged Identification Rates*

Figure 1 shows lagged identification rates and fitted models for one run of each model type and number of identifications per sampling period. When true models, or overly general models, were fit to the simulated data sets, curves generally fit the data well.

The distribution of models selected using AIC and QAIC is shown in Table 1. With population Types A and C, the theoretically correct model was usually selected. However, when using population Type B there was frequent overfitting, with model 3 being selected rather than model 2. There was also some overfitting with the true independent binomial data (population Type BX), but it was less pronounced than for population Type B. When the incorrect model was selected for runs of population Types A and C, then ΔAIC, the difference between the AIC



**Figure 1.** Lagged identification rates ('o', with time lags aggregated so that each time lag interval contained at least 5% of the data points as indicated by the denominator of Eq. (3)) and fitted models of one run for each model type ('A' closed; 'B' permanent emigration; 'C' immigration/emigration from a larger population) and number of identifications per sampling period (10, 20, 40). Fitted models were: 1: – : $R(\tau) = \alpha$; 2: - - : $R(\tau) = \alpha \cdot e^{-\beta \cdot \tau}$; 3: . . . : $R(\tau) = \gamma \cdot e^{-\beta \cdot \tau} + \alpha$.

**Table 1**

Summary of results of model selections on simulated data sets for lagged
identification rates: for each population type and number of identifications per
sampling period, the table gives the number of runs (out of 100) in which models
1, 2, or 3 were selected using either the AIC or QAIC criteria. The correct
selections are shown in bold. The final column tabulates the number of runs in
which QAIC was used (i.e., $\hat{c} > 1$)

| Population type | Identifications per sample | Model selected using AIC | | | Model selected using QAIC | | | QAIC preferred ($\hat{c} > 1$) |
|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 1 | 2 | 3 | |
| A | 10 | **87** | 10 | 3 | **87** | 10 | 3 | 26/100 |
| A | 20 | **87** | 12 | 1 | **87** | 12 | 1 | 0/100 |
| A | 40 | **97** | 3 | 0 | **97** | 3 | 0 | 0/100 |
| B | 10 | 0 | **69** | 31 | 0 | **76** | 24 | 53/100 |
| B | 20 | 0 | **47** | 53 | 0 | **48** | 52 | 29/100 |
| B | 40 | 0 | **34** | 66 | 0 | **35** | 65 | 29/100 |
| C | 10 | 0 | 2 | **98** | 0 | 2 | **98** | 64/100 |
| C | 40 | 0 | 0 | **100** | 0 | 0 | **100** | 45/100 |
| C | 20 | 0 | 0 | **100** | 0 | 0 | **100** | 13/100 |
| BX | 10 | 0 | **79** | 21 | 0 | **88** | 12 | 53/100 |
| BX | 20 | 0 | **76** | 24 | 0 | **79** | 21 | 49/100 |
| BX | 40 | 0 | **84** | 16 | 0 | **87** | 13 | 48/100 |

of the theoretically-correct model and that of the chosen model, was almost always
less than 4.0 indicating substantial support for the theoretically-correct model in
situations when AIC is theoretically valid (Burnham and Anderson, 2002). However,
with population Type B, there was little support for the theoretically correct model
(indicated by $\Delta$AIC or $\Delta$QAIC $> 4$) in 22–24% of the 300 runs.

QAIC was selected over AIC (as $\hat{c} > 1$) in 259 of the 900 runs, especially when
there were fewer data (Table 1). The two criteria selected different models in only
nine of these runs, with QAIC selecting the theoretically correct model in all these
cases. QAIC was also frequently selected with the true independent data (population
Type BX), indicating that the selection of QAIC over AIC does not necessarily
indicate non independent data (Table 1).

Results with $N = 1,000$ (and so $Z = 3,000$) were qualitatively similar in most
respects to those with $N = 100$, and so are not presented in detail. The theoretically
correct model was generally chosen by both AIC (73% of time) and QAIC (76%
of time). When $N = 1,000$, population Type B was overfit less often (12–29%
depending on number of identifications per sample and use of AIC or QAIC)
than with $N = 100$ (24–66%). However, in contrast to runs with $N = 100$, when
$N = 1,000$ population Type C was frequently underfit when there were 10 or
20 identifications per sample (AIC choosing models 1 or 2 in 58% and 14%
of runs, respectively), but not with 40 identifications per sample in which the
theoretically correct model 3 was always chosen. This underfitting may be explained
by the sparsity of repeat identifications of the same individual, which are needed
to discriminate the more complex models; with $Z = 3,000$, the expected number

of individuals identified three or more times is about 0.2 with 10 identifications per sample, and 1.8 with 20 identifications per sample (calculated using binomial probabilities).

Using the second-order information criterion, which corrects for small sample size, made very little difference to the results of the model selection. For population Type B, $AIC_c$ selected model 2 rather than the AIC-selected model 3 for three runs, and $QAIC_c$ selected model 2 rather than the QAIC-selected model 3 for two runs. For all other runs, including all those investigating lagged association rates, the second-order criterion selected the same model as AIC and QAIC.

### 3.2. *Model-Fitting: Lagged Association Rates*

The lagged association rates and fitted models for one run of each model type and number of identifications per sampling period shown in Fig. 2 suggest that true, or overly general, models generally fit the data well. QAIC was almost always chosen over AIC when fitting lagged association rates, indicating overdispersion, and the use of QAIC greatly improved performance in model selection for lagged association rates (Table 2). However, even when using QAIC, with all population



**Figure 2.** Lagged association rates ('o', with time lags aggregated so that each time lag interval contained at least 5% of the data points as indicated by the denominator of Eq. (4)) and fitted models of one run for each model type ('D' random; 'E' casual acquaintances; 'F' permanent companions plus casual acquaintances) and number of groups identified per sampling period (1, 2, 4). Fitted models were: 1: – : $g(\tau) = \alpha$; 4: - - : $g(\tau) = (1 - \alpha) \cdot e^{-\beta \cdot \tau} + \alpha$; 5: . . . : $g(\tau) = ((1 - \alpha) \cdot e^{-\beta \cdot \tau} + \alpha) \cdot e^{-\gamma \cdot \tau}$.

**Table 2**

Summary of results of model selections on simulated data sets for lagged
association rates: for each population type and number of identifications per
sampling period, the table gives the number of runs (out of 100) in which models
1, 4, or 5 were selected using either the AIC or QAIC criteria. The correct
selections are shown in bold. The final column tabulates the number of runs in
which QAIC was used (i.e., $\hat{c} > 1$)

| Population type | Groups identified per sample | Model selected using AIC | | | Model selected using QAIC | | | QAIC preferred ($\hat{c} > 1$) |
|---|---|---|---|---|---|---|---|---|
| | | 1 | 4 | 5 | 1 | 4 | 5 | |
| D | 1 | **61** | 5 | 34 | **81** | 4 | 15 | 100/100 |
| D | 2 | **56** | 8 | 36 | **84** | 3 | 13 | 100/100 |
| D | 4 | **52** | 11 | 37 | **76** | 8 | 16 | 100/100 |
| E | 1 | 0 | **38** | 62 | 0 | **78** | 22 | 99/100 |
| E | 2 | 0 | **28** | 72 | 0 | **54** | 46 | 100/100 |
| E | 4 | 0 | **21** | 79 | 0 | **41** | 59 | 100/100 |
| F | 1 | 0 | **27** | 73 | 1 | **63** | 36 | 100/100 |
| F | 2 | 0 | **11** | 89 | 0 | **33** | 67 | 100/100 |
| F | 4 | 0 | **4** | 96 | 0 | **16** | 84 | 100/100 |
| EX | 1 | 0 | **88** | 12 | 0 | **87** | 13 | 10/100 |
| EX | 2 | 0 | **81** | 19 | 0 | **79** | 21 | 9/100 |
| EX | 4 | 0 | **84** | 14 | 0 | **84** | 16 | 9/100 |

types and sampling rates, there was quite frequent overfitting with overly complex
models being chosen (22–79% of runs as opposed to 12–21% with the independent
binomial data of model EX), but almost no underfitting (Table 2). Quite frequently
(20% of runs), there was little support for the theoretically correct model (indicated
by $\Delta$QAIC$> 4$; Burnham and Anderson, 2002).

### 3.3. *Confidence Intervals of Parameter Estimates*

Table 3 lists the results of the investigation of the width of estimated confidence
intervals produced by different estimation procedures (parameter estimates showed
little evidence of bias). For the lagged identification rates, all the methods give
confidence spans not too far from the true value for all estimated parameters
and similar coverage probabilities (means 84–92%), although the bootstrap method
appeared best in terms of providing confidence spans generally closest to the true
values, and with least variation, as well as the best coverage (mean 92%), close to
the nominal 95%.

For the lagged association rates, the results are less encouraging (Table 3).
Bootstrap estimates are invalid, and likelihood and quasi-likelihood confidence
intervals are generally much too narrow, and have a low probability of containing
the true parameter value. Confidence intervals for the lagged association rates from
the jackknife procedure have widths closest to the true widths, although these also
tend to be somewhat too narrow, to have high variance, and the coverage is less
good than for the quasi-likelihood method (mean coverage probability 48%).

**Table 3**

Comparison of width of confidence intervals of models of lagged identification and association rates. The approximate true width is the range between the 2.5 and 97.5% percentiles of the estimated parameters for 100 random runs. Also shown are the median standard deviations of the widths (in round parentheses) and coverages [in square brackets] estimated over runs using bootstrap (1,000 bootstrap replicates), jackknife, likelihood, and quasi-likelihood methods

| | | | | Median upper 95% c.i. – Lower 95% c.i. (SD): | | | |
|---|---|---|---|---|---|---|---|
| Population | Model | Parameter | True | Bootstrap | Jackknife | Likelihood | Quasi-likelihood |
| *Lagged identification rates* | | | | | | | |
| A | 1 | $\alpha$ | 0.00083 | 0.00083 (0.00012) [98%] | 0.00133 (0.00049) [96%] | 0.00115 (0.00001) [100%] | 0.00115 (0.00001) [100%] |
| B | 2 | $\alpha$ | 0.00535 | 0.00455 (0.00110) [76%] | 0.00579 (0.00268) [85%] | 0.00247 (0.00062) [55%] | 0.00253 (0.00065) [57%] |
| B | 2 | $\beta$ | 0.00296 | 0.00248 (0.00035) [92%] | 0.00263 (0.00105) [77%] | 0.00228 (0.00017) [89%] | 0.00236 (0.00025) [93%] |
| C | 3 | $\alpha$ | 0.31382 | 0.32775 (0.16077) [94%] | 0.40874 (0.26786) [74%] | 0.38301 (0.17513) [87%] | 0.40484 (0.18191) [87%] |
| C | 3 | $\beta$ | 0.00456 | 0.00531 (0.00200) [98%] | 0.00534 (0.00304) [77%] | 0.00652 (0.00249) [94%] | 0.00691 (0.00268) [95%] |
| C | 3 | $\gamma$ | 0.00101 | 0.00092 (0.00016) [91%] | 0.00154 (0.00058) [95%] | 0.00076 (0.00007) [89%] | 0.00080 (0.00008) [91%] |
| *Lagged association rates* | | | | | | | |
| D | 1 | $\alpha$ | 0.02671 | | 0.01935 (0.00735) [39%] | 0.01093 (0.00068) [53%] | 0.02252 (0.00201) [88%] |
| E | 4 | $\alpha$ | 0.06124 | | 0.03133 (0.01992) [43%] | 0.01547 (0.00124) [38%] | 0.04049 (0.00521) [86%] |
| E | 4 | $\beta$ | 0.01345 | | 0.01257 (0.00582) [53%] | 0.00197 (0.00055) [21%] | 0.00518 (0.00158) [55%] |
| F | 4 | $\alpha$ | 0.26979 | | 0.15419 (0.14848) [50%] | 0.02101 (0.00493) [9%] | 0.07664 (0.01931) [44%] |
| F | 4 | $\beta$ | 0.05384 | | 0.03237 (0.17389) [55%] | 0.00239 (0.00816) [11%] | 0.00949 (0.03127) [35%] |

## 4. Discussion

The analyses presented here confirm that maximizing the sum of non independent log-likelihoods produces models which provide good visual fits to lagged identification and lagged association rate data (Figs. 1–2). The results (Table 3) also corroborate the use of the bootstrap method for estimating the precision of parameters of models of lagged identification rates (Whitehead, 2001). As suggested

earlier (Whitehead, 1995), the temporal jackknife, in which data from particular time intervals are omitted sequentially in the calculation of pseudovalues, is a feasible method of estimating the precision of parameters in models of lagged association rates. However, these jackknife estimates of precision are themselves not very precise, may overestimate precision, and quasi-likelihood confidence interval estimates seem to have rather better coverage probabilities (Table 3).

The results of the investigations of using AIC or QAIC for selecting models of lagged identification or lagged association rates are more equivocal, and depend very much on the data set and models being compared. With models of lagged identification rates, the use of QAIC only marginally improved performance over AIC (Table 1), whereas for lagged association rates QAIC generally performed much better (Table 2). In some cases, the criteria performed very well, almost invariably selecting the correct model (Tables 1 and 2). However, with population Types B and F, incorrect models were selected more frequently than correct ones, and sometimes there was very little support for the correct model. Except in the case of very sparse data (1,000 individuals and a low sampling rate), overfitting was much more of a problem than underfitting, as found in other simulation studies of the effectiveness of AIC (e.g., Andres and Currim, 2003), and predicted by the theoretical work of Woodroofe (1982). Overfitting was also present with true independent binomial data (models BX and EX; Tables 1 and 2), but less prevalent than with realistic data.

Although the author used several different types of data and sampling rates, as well as some of the most obvious models that can be fit to lagged identification and lagged association rate data, he only examined a very small portion of the types of data and models that are possible. Furthermore, he has not used simulation to examine either lagged identification rates between areas (Whitehead, 2001), or standardized lagged association rates. Given the large variation in the success of model selection using AIC and QAIC in this sample, we conclude that these criteria should be used cautiously with lagged identification and lagged association rate data. However, even with true independent and well-distributed data, these criteria are not always successful (Tables 1 and 2; Burnham and Anderson, 2002), so the author does think that, when used cautiously, AIC or QAIC may be useful in informing us about movements and associations of identified individuals through lagged identification and association rates.

Both lagged identification and association rates are principally descriptive measures, they illustrate how residence within areas and associations between animals change with time. As is the case with population Types E and F, quite different processes can produce the same pattern of lagged association, or identification, rates. Thus, the patterns observed, and models fit, do not prescribe the underlying process. There are methods for fitting mechanistic models to some types of movement data (e.g., Hilborn, 1990), but, to my knowledge, the corresponding techniques have yet to be developed for social analyses.

## Appendix

### *Theoretical Parameter Estimates for Lagged Association Rates*

In models 1 and 4, $\alpha$ is the probability of being in the same group again after a very long time, which for populations D and E is simply $1/\text{No. groups} = 1/W$.

The lagged association rate, $g(\tau)$, is the probability that if $X$ and $Y$ are identified in the same group, then, $\tau$ time units later, $Y$ is once again in the same group as $X$. Then, in the independent switching model of population Type E, given $g(\tau)$:

$$g(\tau + \delta\tau) = g(\tau) \cdot [(1 - (\pi\delta\tau))^2 + (\pi\delta\tau)^2 \cdot s]$$
$$+ (1 - g(\tau)) \cdot [2 \cdot \pi\delta\tau \cdot (1 - \pi\delta\tau) \cdot s + (\pi\delta\tau)^2 \cdot s]$$

where $\pi\delta\tau$ is the probability that an individual switches to a different group in $\delta\tau$, and $s$ is the probability that, given the individual switches groups, it switches to a particular group (containing the other individual). Ignoring the small terms in $(\delta\tau)^2$:

$$\delta g(\tau) = -[g(\tau) - (1 - g(\tau)) \cdot s] \cdot 2\pi\delta\tau.$$

Integrating, and using $g(0) = 1$, gives:

$$g(\tau) = 1/(1 + s)e^{-2\cdot(1+s)\cdot\pi\cdot\tau} + s/(1 + s).$$

In population Type E, $\pi = \lambda$ and $s = 1/(W - 1)$. Thus, with model 4, $\alpha = 1/W$ and $\beta = 2 \cdot \lambda \cdot W/(W - 1)$.

In population Type F with permanent units of size $U$, let the probability that grouped individuals are also units members be $u$. Then:

$$g(\tau) = u + (1 - u) \cdot (1/(1 + s)e^{-2\cdot(1+s)\cdot\pi\cdot\tau} + s/(1 + s)).$$

Now, the probability that two individuals are members of the same unit is $1/U$, and the probability that they are members of different units but the same group is $(1 - 1/U)/W$. Thus, $u = W/(W + U - 1)$, and $\pi$ and $s$ are as for population Type E. Substituting:

$$g(\tau) = [W^2 + U - 1]/[W \cdot (W + U - 1)]$$
$$+ (U \cdot W - W - U + 1)/[W \cdot (W + U - 1)] \cdot e^{-2\cdot\lambda\cdot(W/(W-1))\cdot\tau}.$$

From this, for model 4, $\alpha = (W^2 + U - 1)/[W \cdot (W + U - 1)]$ and $\beta = 2 \cdot \lambda \cdot W/(W - 1)$.

## Acknowledgments

## References

Akaike, H. (1973). Information theory as an extension of the maximum likelihood principle. In: Petrov, B. N., Csaki, F., eds. *Second International Symposium on Information Theory*. Budapest, Hungary: Akademiai Kiado.

Andres, R. L., Currim, I. S. (2003). Retention of latent segments in regression-based marketing models. *Int. J. Res. Market.* 20:315–321.

Bejder, L., Fletcher, D., Bräger, S. (1998). A method for testing association patterns of social animals. *Animal Behav.* 56:719–725.

Burnham, K. P., Anderson, D. R. (2002). *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach.* New York: Springer-Verlag.

Edwards, A. W. F. (1992). *Likelihood.* Baltimore, MD: John Hopkins University Press.

Efron, B., Stein, C. (1981). The jackknife estimate of variance. *Ann. Statist.* 9:586–596.

Gowans, S., Whitehead, H., Arch, J. K., Hooker, S. K. (2000). Population size and residency patterns of northern bottlenose whales (*Hyperoodon ampullatus*) using the Gully, Nova Scotia. *J. Cetacean Res. Manage.* 2:201–210.

Hilborn, R. (1990). Determination of fish movement patterns from tag recoveries using maximum likelihood estimators. *Canad. J. Fisheries Aquatic Sci.* 47:635–643.

Karczmarski, L., Würsig, B., Gailey, G., Larson, K. W., Vanderlip, C. (2005). Spinner dolphins in a remote Hawaiian atoll: social grouping and population structure. *Behav. Ecol.* 16:675–685.

Ottensmeyer, C. A., Whitehead, H. (2003). Behavioural evidence for social units in long-finned pilot whales. *Canad. J. Zool.* 81:1327–1338.

Pledger, S., Pollock, K. H., Norris, J. L. (2003). Open capture-recapture models with heterogeneity: I. Cormack-Jolly-Seber model. *Biometrics* 59:786–794.

White, G. C., Burnham, K. P. (1999). Program MARK: survival estimation from populations of marked animals. *Bird Study* (Supplement) 46:120–138.

Whitehead, H. (1995). Investigating structure and temporal scale in social organizations using identified individuals. *Behav. Ecol.* 6:199–208.

Whitehead, H. (2001). Analysis of animal movement using opportunistic individual-identifications: application to sperm whales. *Ecology* 82:1417–1432.

Woodroofe, M. (1982). On model selection and the arc sine laws. *Ann. Statist.* 10:1182–1194.