

BU-911-MB

Maximum likelihood estimation to assign paternity

within multiply sired broods:

use of the EM algorithm

Charles E. McCulloch
Biometrics Unit
Cornell University
Ithaca, New York 14853

*Janis L. Dickinson
Department of Entomology
Cornell University
Ithaca, New York 14853

* Present address: Department of Zoology, Arizona State University,
Tempe, Arizona 85287

Abstract – Attempts to measure male reproductive success in natural populations have been hindered by the fact that paternity often cannot be deduced from behavioral data alone. Recently, there has been increasing reliance on the use of molecular polymorphism to infer genealogical relationships. Such inference usually requires the use of statistical procedures to resolve ambiguities. We advocate the use of the EM algorithm (Dempster et al. 1977) to calculate maximum likelihood estimators of probabilities of parentage for each of a set of suspected sires. This method permits the researcher to apportion a multiply sired brood among a female's different mates. For this purpose, the maximum likelihood method is better than the LOD ratio method proposed by Meagher (1986). Simulations showed that the estimators are usually quite accurate with brood sizes of 25 or greater and that the probability that the male with the highest paternity will be ranked first is high with brood sizes as low as 10.

In recent years, there has been increasing use of electrophoretic (Meagher 1986) and other codominant variants (Jeffreys et al., 1985a; 1985b) to infer genealogical relationships in populations of animals and plants. Multiple parentage within broods has been found to be relatively widespread (Hanken and Sherman, 1981; McCauley and O'Donnell, 1984; Ellstrand and Marshall, 1986). As maternal parents are more easily identified than are male parents, it is usually the paternity of offspring that is in question. When a brood may be sired by more than one male, paternity can be inferred using data on genetic markers in conjunction with statistical estimation procedures based upon Mendelian transmission probabilities. Because, in most cases, offspring possess putative genotypes that could have been derived from more than one of a set of males, it is rare that paternity can be unambiguously determined, even in laboratory studies (Dickinson 1986, Foltz and Pashley 1986). Ambiguity comes about in three ways: (1) Males may be identical at each polymorphic locus examined; if this is the case, the genetic markers used are not informative. (2) Males may be identical for at least one of two putative alleles at each locus. (3) Even if all of the males are different at a locus, paternity may still be ambiguous if the female is heterozygous and shares an allele in common with each of two suspects. Because ambiguities are common, accurate quantitative estimates of parentage are needed.

Statistical likelihood has been employed to assign paternity in studies involving humans (Walker, 1983; Thompson, 1986), other animals (Foltz and Hoogland, 1981; Dickinson, 1986; Foltz and Pashley, 1986), and plants (Meagher, 1986; Schoen and Stewart, 1986). In cases in which a brood may be sired by more than one male, there have been two basic approaches. The first and

most widely utilized approach involves use of log-likelihood (LOD) ratios (Meagher 1986; Meagher and Thompson, 1986). Use of LOD ratios takes into consideration the likelihood that male "a" fathered a particular offspring with a particular female in relation to the likelihood that he is a randomly sampled member of the parental generation (Meagher, 1986). LOD ratios are compared among potential male parents and the male with the highest ratio is assigned paternity for the offspring in question. Each offspring is considered independently, such that the number of offspring used to acquire each estimate is one. The second method is maximum likelihood (ML) estimation based upon the Mendelian transmission probabilities of all potential sires (suspects) and the relative proportions of different offspring genotypes in the brood (Dickinson, 1986; Foltz and Pashley, 1986). The clutch or brood is then assigned to males in proportions corresponding to the males' probabilities of paternity, and the number of offspring used to calculate each estimate is equal to the total number of offspring in the brood.

In this paper we demonstrate how the ML estimation procedure described in Dickinson (1986) and Foltz and Pashley (1986) can be extended to include cases in which there are more than two potential sires, and demonstrate how the EM-algorithm can be used to find ML estimators. We also evaluate the effects of sample size on accuracy and bias of the ML estimators and compare ML estimation with least squares estimation procedures.

MULTIPLE PATERNITY WHEN MATERNITY IS KNOWN

We will examine the case in which there is a single mother and the genotypes of the mother, offspring, and male suspects are known. Maternity is usually certain in laboratory studies of sperm utilization patterns in insects (McCauley and Reilly, 1984; Dickinson, 1986; Foltz and Pashley, 1986), as well as in field studies of plants (Meagher, 1986; Schoen and Stewart, 1986), most species of animals that bear live young (Hanken and Sherman, 1981), and species that nest or exhibit maternal care of eggs (but see Yom-Tov et al., 1974; Gowaty and Karlin, 1986; Tallamy, 1986). The probabilities that we will examine are the probabilities of paternity for each male suspected of having mated with a given female. Although the analysis is dependent on the female's genotype, the fact that she remains constant throughout allows us to use simplified notation that does not outwardly reflect this dependence.

Our objective is to come up with an estimate of the proportion of offspring sired by each of a set of male suspects. In the case of laboratory studies, these suspects include all of a female's sequential mates. When animals are observed in the field, it may be rare that copulations are actually witnessed (Gavin and Bollinger, 1985; Mumme et al., 1985). In such cases, suspects may include the males whose territories overlap that of the female in addition to other males she has associated with during the study.

Let us first establish some notation:

S = number of "suspects"

G = number of distinct genotypes among the offspring

$P^*(j|i)$ = probability of an offspring having the j th genotype, assuming the i th male was the sire

$P^*(j)$ = probability of an offspring being the j th genotype

$P(i)$ = probability that the i th male will sire an offspring

$P(i|j)$ = probability of the i th male being the sire of an offspring of genotype j

$f(j)$ = number of offspring of the j th genotype in the sample

As shown by Schoen and Stewart (1986) and, for two males, by Foltz and Pashley (1986), $P(j)$ is a linear function of the $P(i)$'s:

$$P^*(j) = P^*(j|1) P(1) + P^*(j|2) P(2) + \dots + P^*(j|S) P(S). \quad (1)$$

In equation (1), $P^*(j)$ can be calculated and the statistical problems center on estimating the $P(i)$'s from the data.

THE USE OF MAXIMUM LIKELIHOOD ESTIMATION

Thompson (1986, Appendix 3) provides a general description of maximum likelihood equations and details their use in genealogy reconstruction. Description of a maximum likelihood estimation procedure that is most similar to ours can be found in Foltz and Pashley (1986) for the case in which there are two

potential sires. For the problem we consider, the likelihood is given by

$$L(P(1), \dots, P(S)) = P^*(1)^{f(1)} P^*(2)^{f(2)} \dots P^*(G)^{f(G)}, \quad (2)$$

and the objective for any given set of data is to find the values for the $P(i)$'s that maximize it.

This can be easily accomplished for two males since the probability of the first male siring offspring is just one minus the probability of the second. The likelihood can be calculated for a grid of finely spaced values (say, every .001) for the probability of the first male using a simple program, a spreadsheet package, or a statistics package like MINITAB. This can also be plotted to view the entire likelihood. The value that gives the largest value for the likelihood is the maximum likelihood estimator.

If there are more than two males, the problem is more complicated. The likelihood can still be evaluated for a grid of finely spaced values, but this quickly becomes time-consuming. There are many numerical algorithms available to maximize nonlinear functions such as (2), for example, the Newton-Raphson technique (Kennedy and Gentle, 1980, p. 442). However, a simpler method that works very reliably for this problem is the EM algorithm (Dempster et al., 1977). The EM algorithm starts with a guess as to the values of the $P(i)$ and iteratively calculates new values that increase the value of the likelihood. Iterations continue until the estimates fail to change and the likelihood is no longer increased. It works as follows:

0. Obtain initial estimates of $P(1), P(2), \dots, P(S)$.
1. Using the current estimates of $P(1), P(2), \dots, P(S)$, calculate the $P(i|j)$ using Bayes' formula (see step 1 below).
2. Use the $P(i|j)$ to apportion $f(j)$ into estimated frequencies attributable to each "suspect".
3. Sum the frequencies for the i th male over all of the genotypes and use the sum to get a new estimate of $P(i)$.
4. Continue to iterate steps 1 through 3 until successive estimates of $P(1), P(2), \dots, P(S)$ change very little.

More specifically, if we denote the portion of the $f(j)$ apportioned to the i th male as $f(i,j)$ and if estimates at the k th iteration are denoted by a superscript k , in brackets, the algorithm is:

0. $k=0, P^{(k)}(i) = 1/S$.
1. $k=k+1, P^{(k)}(i|j) = \frac{P^*(j|i)P^{(k-1)}(i)}{\sum_r P^*(j|r)P^{(k-1)}(r)}$
2. $f(i,j)^{(k)} = f(j)P^{(k)}(i|j)$
3. $P^{(k)}(i) = \sum_j f^{(k)}(j|i) / \sum_j f(j)$
4. If $\max_i \{|P^{(k)}(i) - P^{(k-1)}(i)|\} > \text{tolerance value}$, return to step 1, otherwise stop.

This algorithm is easily programmed and a BASIC program to do the computations

that runs on an IBM PC is available from the authors. Table 1 illustrates how the algorithm is implemented.

Several comments are in order to describe the performance of this algorithm in more detail:

1. Estimates of the $P(i)$ and the $P^*(j)$ will always be between zero and one (as they should be). This is not true of other estimation techniques, such as least squares.
2. If the algorithm is used for data which unambiguously indicate that all of the offspring come from one male, then the probability of that male siring offspring is estimated to be one and the rest are estimated to be zero.
3. As an extension of the case in 1, if all of the data are unambiguous and more than one male sires offspring within the brood, the probabilities $[P(i)]$ for each male are estimated to be the sample relative frequencies that may be unambiguously attributed to each male.
4. In cases in which genetic patterns for two or more males are identical, the data give no information for distinguishing among them. The way the algorithm is implemented (starting with equal probabilities for each male), the probabilities $[P(i)]$ estimated for those males will be identical.

ASSUMPTIONS

There are a number of assumptions, both genetic and statistical, inherent in the proper use and interpretation of this model. As written, the model assumes that multiple maternity will not occur. The model can be easily adapted to handle the case of multiple maternity with single paternity by interchanging the roles of males

and females. It can also be adapted to handle the case of multiple maternity *and* paternity by considering every male-female pair, but $P^*(j|i)$, $P(i)$, $P^*(j)$, and $P(i|j)$ would have to be redefined as follows:

$P^*(j|i)$ = probability of an offspring having the j th genotype assuming it is an offspring of the i th male-female pair.

$P^*(j)$ = probability of an offspring being the j th genotype.

$P(i)$ = probability that the i th male-female pair will produce an offspring.

$P(i|j)$ = probability of the i th male-female pair producing an offspring of genotype j .

Simultaneous analysis of multiple maternity and paternity would require a large amount of data since each pairing of a male and female would introduce a parameter $[P(i)]$ to be estimated.

In order to calculate the probabilities of genotypes $[P^*(j/i)]$, we will need to make the following assumptions:

- i. Mendelian inheritance
- ii. No correlations among loci.

However, if loci are linked and the joint multilocus probabilities can be calculated, then this method can still be used. If population frequencies are used for $P^*(j|i)$ rather than the true value for individuals, then random mating must be assumed.

Offspring are treated as statistically independent in the formation of the likelihood equation (2). This is likely to be a good assumption under conditions in

which multiple mating is known to have occurred. It may be a poorer assumption when mating is rarely observed and one male is likely to sire the entire litter. If litters are rarely multiply sired, but the identity of the sire changes from one litter to the next, then the assumption of independence will not be valid. The model would need to be rewritten with litter, rather than single offspring, as the unit of observation. The same general approach could then be used providing that data on a sufficient number of litters are available.

The model assumes that all potential sires can be identified and that their genetic patterns are known. This might be a problem with field data. In cases in which copulations are difficult to observe, the risk of leaving a critical male out of the analysis is relatively high. The problem of which males to include in the analysis is one that plagues all field studies of multiple mating and paternity.

PERFORMANCE OF THE MAXIMUM LIKELIHOOD ESTIMATORS

We first compare the ML method with the LOD method of Meagher (1986). Meagher's method is an example of a "classification" technique, in which each offspring is assigned unambiguously to a sire. In the case of a single mother, this method corresponds to assigning each offspring to the sire with the largest probability of siring an offspring of that genotype. Classification techniques have been found to perform very poorly (Bryant and Williamson, 1978). A simple example serves to illustrate the problem. Suppose we have a situation with two males and two genotypes with the $P^*(j/i)$ given in Table 2. Using the LOD ratio method, all offspring of genotype 1 will be assigned to male 1 and all offspring of

genotype 2 will be assigned to male 2. Thus the estimated proportion of offspring attributable to male 1 will be the proportion of offspring of genotype 1. This is not estimating the true proportion of offspring attributable to male 1 [$P(1)$], but instead is estimating the probability of genotype 1, which is equal to $(0.75) P(1) + (0.5) P(2) = 0.5 + (0.25) P(1)$, since $P(2) = 1 - P(1)$. Even with arbitrarily large numbers of offspring per female, use of the LOD ratio will give unreliable estimates. The ML method, on the other hand, converges to the true value as the sample size increases. Hence, Meagher's method cannot be recommended when the goal is to use all of the information in the sample. In essence, Meagher's method uses samples of size 1 (each offspring is considered separately) and any inaccuracies in the likelihood method due to small sample sizes are perpetuated as the number of offspring per female increases.

We next investigate the accuracy of the ML estimators and compare them with the more easily obtained least squares estimators. Least squares estimators have the advantage of being unbiased (their average value in replicated experiments is the true value) while maximum likelihood estimators are not. To compare them, we must therefore consider both bias and variance. A common measure of closeness of the estimators to the true value is the mean square error (average squared difference between estimates and true value). We have used this measure to evaluate the ML estimators and to compare them to the least squares estimators. Twelve separate sets of simulations were performed to evaluate the estimators (APPENDIX 1). One set of simulations was chosen to reflect the range of gene patterns found in Dickinson (1986).

Figures 1 through 3 display the performance of the ML estimators for parameter configuration D (APPENDIX 1). This configuration was neither the best nor the worst case for performance of the ML estimators. Figure 1 shows the bias, which is very small in estimating any of the parameters, even for small sample sizes. The worst case is estimating $P(1)$ with samples of size 4. In this case the ML estimator yields a $P(1)$ that is too low by 0.036 on average (the true value is 0.5, while the average of the estimator is 0.474). Cases where the bias was not small were typically cases where the true probability was close to zero or one. To understand this, consider the case in which the true probability is close to one. Often, the ML estimator will be equal to or close to one. However, whenever it is not, it will be less than one and hence the mean of the estimator will also be less than one. The least squares estimator balances those values less than one with some values greater than one, so that they average out to one. In this case, unbiasedness requires zero variance or values out of the range of zero to one. Hence, unbiasedness is probably not a good property to require. These considerations are demonstrated in Table 3.

The estimates of mean square error from the simulations showed that with small sample sizes the ML estimators are not terribly accurate. Mean square error is equal to the sum of the square of the bias and the variance. The largest portion of this was usually the variance. Figure 2 shows the relationship between the standard deviation and the sample size. It was not until sample sizes of 25 or so were reached that the variances came down to acceptable levels. This suggests that either large litter sizes or a large number of litters would be necessary to accurately obtain paternity estimates.

We also investigated the probability that the male with a higher probability of paternity would be estimated to have a higher paternity value. Fig. 3 shows that a reasonably high probability of correct ranking ($\Pr\{CR\}$) can be achieved, even with fairly small sample sizes.

Unfortunately, the "usual" approximate tests and confidence intervals for maximum likelihood estimators are not valid for this model because of the frequency of maximum likelihood estimates that are exactly zero or one. This means, for example, that the chi-square tests of the hypothesis that $P(1)$ equals zero or one recommended in Foltz and Pashley (1986) are invalid. It also means that the "usual" methods of calculating standard errors for maximum likelihood estimators may give misleading results.

Because researchers usually do not know in advance whether multiple paternity occurs, it was important to investigate the performance of the estimator when one male sired all of of a litter (i.e. the data were not independent), but different litters could be sired by different males. In doing so, we found that the true male had the highest estimated $P(i)$, in almost all cases, and that the $P(i)$ was frequently estimated to be unity. For example, in configuration D (APPENDIX 1) with only 4 offspring, the correct male had the highest $P(i)$ in about 70% of the cases; with 50 offspring the percentage increased to 90%. This suggests that the estimators can be fairly reliable for giving rank order information when there is a lack of independence, even for small sample sizes.

CONCLUSIONS

Assigning paternity on the basis of isozyme variants and other molecular polymorphisms is problematic, even in cases in which behavioral data on suspects are good (McCracken and Wilkinson, in press). The ML estimation procedure we have described is best applied to organisms with large clutch or litter sizes, such as certain fish (Darling et al., 1980), reptiles (Gibson and Falls, 1975), amphibians (Tilley and Hausman, 1976), mollusks (Murray, 1964), and arthropods (Sassaman, 1978; McCauley and Reilly, 1984). Its usefulness in assigning parentage for many species of mammals (Hanken and Sherman, 1981) and birds (Gowaty and Karlin, 1984; Mumme et al., 1985) will be limited; we don't recommend its use for species with small numbers of offspring (fewer than 10) unless data are available for a large number of families.

In cases in which the numbers of offspring are sufficient to justify use of this method, researchers will be able to answer questions about parentage with fewer families than were previously needed. For example, Dickinson (1986, in press) used ML estimation in conjunction with isozyme data to quantify the proportion of offspring fathered by each of a female's two consecutive mates in laboratory studies of determinants of paternity in the milkweed leaf beetle. The resulting ML estimators for second males were compared among treatments in which mating duration varied using a Mann-Whitney U-test. In situations like this, ML can be employed to ask questions about the behavioral and morphological determinants of male reproductive success. We propose use of ML estimation as an alternative to the LOD ratio method (Meagher and Thompson 1986) in situations in which it is likely

BU-911-MB

that there is multiple paternity within clutches. Unlike the LOD ratio, the ML estimator we describe does not usually permit the researcher to infer that a particular father sired a particular offspring. However, the fact that it is based upon the relative numbers of offspring of different genotypes in the sample makes it a more desirable method for dividing a multiply sired brood among a female's different mates.

ACKNOWLEDGMENTS

We would like to thank G. Eickwort, W. Koenig, P. Kukuk, W. Sheehan, and P. Sherman for reading an early draft of this manuscript. The paper was substantially revised according to comments from two anonymous reviewers; we are especially grateful for their suggestions and criticisms. Financial support for the work came from the Biometrics Unit and Department of Entomology at Cornell University.

APPENDIX 1

All simulations were run on an IBM PC-AT. The simulation programs were written in the matrix language GAUSS. Random number generation was performed using built-in GAUSS functions, which use a multiplicative congruential generator. Computation of the maximum likelihood estimators is described in the text. The least squares estimators are the usual least squares estimates found by regressing the observed genotype frequencies on their means [a linear function of the sire probabilities, $p(i)$]. The sire probability estimates were restricted to sum to one. The parameter configurations are given in the table below.

Simulation set	S	G	# Replications	P*(j/i)			# Offspring (NOBS) and P(i)	
A	3	3	1000*	0.5	0.25	0.25	P(i) = (0.6, 0.35, 0.05)	
				0	0.75	0.25	NOBS = 4, 10, 25, 50, 100	
				0.875	0.125	0		
B	2	3	1000	0.5	0.25	0.25	P(i) = (0.75, 0.25)	
				0.875	0.125	0	NOBS = 4, 10, 25, 50, 100	
C	3	4	1000	0.5	0.25	0.25	0	P(i) = (0.5, 0.3, 0.2)
				0	0.75	0.25	0	NOBS = 4, 10, 25, 50, 100
				0	0.25	0.5	0.25	
D	3	3	1000	Same as A			Same as C	

(Appendix 1, continued)

S	2	2	1000	0.5 1.0	0.5 0				P(i) = (1,0), (.95, .05), (.9, .1), (.8, .2), (.7, .3), (.6, .4), (.5, .5), (.4, .6), (.3, .7), (.2, .8), (.1, .9), (.05, .95), (0, 1); NOBS = 10
T	2	3	1000	0.5 0.25	0.5 0.5	0 0.25			P(i) = same as S; NOBS = 10
U	2	3	1000	0.25 0.5 0	0.25 0.5 0	0.25 0 0	0.25 0 0.5	0 0 0.5	P(i) = (1, 0, 0), (.8, .2, 0), (.8, 0, .2), (.6, .4, 0), (.6, .2, .2), (.6, 0, .4), (.4, .6, 0), (.4, .4, .2), (.4, .2, .4), (.4, 0, .6), (.2, .8, 0), (.2, .6, .2), (.2, .4, .4), (.2, .2, .6), (.2, 0, .8), (0, 1, 0), (0, .8, .2), (0, .6, .4) (0, .4, .6), (0, .2, .8), (0,0,1), (.333, .333, .333) NOBS = 10
V	3	5	1000	0.25 0.5 0	0.25 0.5 0	0.25 0 0	0.25 0 0.5	0 0 0.5	P(i) = same as U; NOBS = 25
W	2	4	1000	0.25 0.5	0.25 0.5	0.25 0	0.25 0		P(i) = same as S; NOBS = 25
X	2	5	1000	0.25 0	0.25 0	0.25 0	0.25 0.5	0 0.5	P(i) = same as S; NOBS = 25
Y	2	3	1000	0.5 0.25	0.5 0.5	0 0.25			P(i) = same as S; NOBS = 25

(Appendix 1, continued)

Z	2	2	1000	0.5	0.5	P(i) = same as S; NOBS = 25
				1	0	

* 250 replications when NOBS = 4.

LITERATURE CITED

- BRYANT, P. AND J. A. WILLIAMSON. 1978. Asymptotic behaviour of classification maximum likelihood estimates. *Biometrika* 65: 273-281.
- DARLING, J. D. S., M. L. NOBLE, AND E. SHAW. 1980. Reproductive strategies in the surfperches. I. Multiple insemination in natural populations of the shiner perch, *Cymatogaster aggregata*. *Evolution* 34: 271-277.
- DEMPSTER, A. P., N. P. LAIRD, AND D. B. RUBIN. 1977. Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Stat. Soc., Series B*, 39: 1-38.
- DICKINSON, J. L. 1986. Prolonged mating in the milkweed leaf beetle, *Labidomera clivicollis clivicollis* (Coleoptera: Chrysomelidae): a test of the 'sperm-loading' hypothesis. *Behav. Ecol. Sociobiol.* 18: 331-338.
- ELLSTRAND, N. C. AND D. L. MARSHALL. 1986. Patterns of multiple paternity in populations of *Raphanus sativus*. *Evolution* 40: 837-842.
- FOLTZ, D. W. AND J. L. HOOGLAND. 1981. Analysis of the mating system in the black-tailed prairie dog (*Cynomys ludovicianus*) by likelihood of paternity. *J. Mammol.* 62: 706-712.
- FOLTZ, D. W. AND D. P. PASHLEY. 1986. Estimating the degree of sperm precedence in laboratory mating experiments: a maximum likelihood method. *J. Hered.* 77: 477-478.
- GAVIN, T. AND E. BOLLINGER. 1985. Multiple paternity in a territorial passerine: the bobolink. *The Auk* 102: 550-555.

- GIBSON, A. R., AND J. B. FALLS. 1975. Evidence for multiple insemination in the common garter snake, *Thamnophis sirtalis*. *Can. J. Zool.* 53: 1362-1368.
- GOWATY, P. A., AND A. A. KARLIN. 1984. Multiple maternity and paternity in single broods of apparently monogamous eastern bluebirds (*Sialis sialis*). *Behav. Ecol. Sociobiol.* 15: 91-95.
- HANKEN, J. AND P. W. SHERMAN. 1981. Multiple paternity in Belding's ground squirrel litters. *Science* 212: 351-353.
- JEFFREYS, A. J., V. WILSON, AND S. L. THEIN. 1985a. Hypervariable 'minisatellite' regions in human DNA. *Nature* 314: 67-73.
- JEFFREYS, A. J., V. WILSON, AND S. L. THEIN. 1985b. Individual-specific 'fingerprints' of human DNA. *Nature* 316: 76-79.
- KENNEDY, W. I. AND J. E. GENTLE. 1980. *Statistical computing*. Marcel Dekker, New York.
- MCCAULEY, D. E. AND R. O'DONNELL. 1984. The effect of multiple mating on genetic relatedness in larval aggregations of the imported willow leaf beetle *Plagioderma versicolora* (Coleoptera: Chrysomelidae). *Behav. Ecol. Sociobiol.* 15: 287-291.
- MCCAULEY, D. E. AND L. K. REILLY. 1984. Sperm storage and sperm precedence in the milkweed beetle, *Tetraopes tetraophthalmus* (Forster) (Coleoptera: Cerambycidae). *Ann. Entomol. Soc. Amer.* 77: 526-530.
- MCCRACKEN, G. F. AND G. S. WILKINSON. (in press) Allozyme techniques and kinship assessment in bats. Kunz, T.H., ed., *Behavioral and Ecological Techniques for Research on Bats*, Smithsonian Press, Washington, D.C.

- MEAGHER, T. R. 1986. Analysis of paternity within a natural population of *Chamaelirium luteum*. 1. Identification of most-likely male parents. *Amer. Nat.* 128: 199-215.
- MEAGHER, T. R. AND E. A. THOMPSON. 1986. The relationship between single parent and parent pair genetic likelihoods in genealogy reconstruction. *Theor. Popul. Biol.* 29: 87-106.
- MUMME, R. L., W. D. KOENIG, R. M. ZINK, AND J. A. MARTEN. 1985. Genetic variation and parentage in a California population of acorn woodpeckers. *The Auk* 102: 305-312.
- MURRAY, J. 1964. Multiple mating and effective population size in *Cepaea nemoralis*. *Evolution* 18: 283-291.
- SASSAMAN, C. 1978. Mating systems in porcellionid isopods: multiple paternity and sperm mixing in *Porcellio scaber*. *Latr. Heredity* 41: 385-397.
- SCHOEN, D. J. AND S. C. STEWART. 1986. Variation in male reproductive investment and male reproductive success in white spruce. *Evolution* 40: 1109-1120.
- TALLAMY, D. W. 1986. Age specificity of 'egg dumping' in *Gargaphia solani* (Hemiptera: Tingidae). *Anim. Behav.* 34: 599-603.
- THOMPSON, E. A. 1986. *Pedigree Analysis in Human Genetics*. Johns Hopkins University Press, Baltimore, MD.
- TILLEY, S. G. AND J. S. HAUSMAN. 1976. Allozymic variation and occurrence of multiple inseminations in populations of the salamander, *Desmognathus ochrophaeus*. *Copeia* 4: 734-741.
- WALKER, R. H., ED. 1983. *Inclusion probabilities in parentage testing*. American Association of Blood Banks, Arlington, VA.

YOM-TOV, Y., G. M. DUNNET, AND A. ANDERSON. 1974. Intraspecific nest parasitism in the starling *Sturnus vulgaris*. Ibis 116: 87-90.

BU-911-MB

TABLE 1. Example of use of the EM algorithm to calculate P(i) for two loci.

Alleles are "F" (fast), "M" (medium), and "S" (slow).

a. *Putative Genotypes of Parents*

Hypothetical Loci	<u>Genotypes of Males</u>				<u>Genotype of Female</u>
1	MF	FF	SM	SF	FF
2	MM	SM	SS	SM	SM

b. *P*(j|i) and Observed frequencies of Offspring Genotypes*

Male	<u>Offspring Genotypes (Locus 1 / Locus 2)</u>								
	FF/MS	FF/MM	MF/SM	MF/MM	FF/SS	MF/SS	SF/MM	SF/SM	SF/SS
P*(j 1)	0.25	0.25	0.25	0.25	0	0	0	0	0
P*(j 2)	0.5	0.25	0	0	0.25	0	0	0	0
P*(j 3)	0	0	0.25	0	0	0.25	0	0.25	0.25
P*(j 4)	0.25	0.125	0	0	0.125	0	0.125	0.25	0.125
# of Observations	3	7	2	9	1	5	6	0	1

c. *EM Algorithm*

Iteration	<u>Estimate of P(i)</u>			
	P(1)	P(2)	P(3)	P(4)
0	0.250	0.250	0.250	0.250
1	0.399	0.146	0.259	0.196
2	0.462	0.087	0.266	0.184
3	0.506	0.036	0.279	0.184
⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮
28	0.526	0.000	0.296	0.178
29	0.526	0.000	0.296	0.178

Convergence Reached

TABLE 2. Probabilities of genotypes for two males [$P^*(j/i)$].

Male	Genotype	
	1	2
1	0.75	0.25
2	0.5	0.5

BU-911-MB

TABLE 3. Example of comparison of the maximum likelihood and least squares estimators ($S = 2, G = 2, P(1) = 1, P(2) = 0, P^*(1|1) = 0.5, P^*(2|1) = 0.5, P^*(1|2) = 0.75, P^*(2|2) = 0.25$). For the case in which male 1 sired all ten offspring within the brood, we calculated the probability of drawing each of ten different offspring genotype combinations $[\text{Pr}\{f(1) \text{ and } f(2)\} \text{ when } f(1) = 0,1,2,\dots,10 \text{ and } f(2) = 10,9,8,\dots,0]$. We then compared the resulting estimates of $P(1)$ obtained by ML estimation with those obtained using least squares estimation procedures.

a. *Example*

Oerved Combinations in the Sample <u>f(1)</u> <u>f(2)</u>	<u>Probability of Occurrence</u>	<u>Value of ML Estimator for P(1)</u>	<u>Value of Least Squares Estimator for P(1)</u>
10 0	0.00098	0	-1.0
9 1	0.00977	0	-0.6
8 2	0.04395	0	-0.2
7 3	0.11719	0.2	0.2
6 4	0.20508	0.6	0.6
5 5	0.24609	1	1
4 6	0.20508	1	1.4
3 7	0.11719	1	1.8
2 8	0.04395	1	2.2
1 9	0.00977	1	2.6
0 10	0.00098	1	3.0

b. *Summary:*

MLE:	Bias= -0.23	Variance= 0.11	Mean Sq. Error= 0.16
Least Squares:	Bias= 0	Variance= 0.40	Mean Sq. Error= 0.40

FIGURE LEGENDS

FIG. 1. The relationship between bias and sample size for simulation D (described in APPENDIX 1).

FIG. 2. The relationship between standard deviation of the ML estimates and sample size for simulation D (described in APPENDIX 1).

FIG. 3. The relationship between probability of correct ranking [Pr(CR)] and sample size for simulation D (described in APPENDIX 1). $\Pr\{1>3\}$ is the probability that the male with the highest probability of paternity is assigned a higher paternity value than the male with the lowest probability of paternity. $\Pr\{1>2\}$ is the probability that the ML estimator for the male with the highest probability of paternity is greater than the ML estimator for the male with the second highest probability of paternity. $\Pr\{2>3\}$ is the probability that the ML estimator for the male with the second highest probability of paternity is greater than the ML estimator for the male with lowest probability of paternity.

Biases of MLEs of P1, P2 and P3

with Increasing Sample Sizes

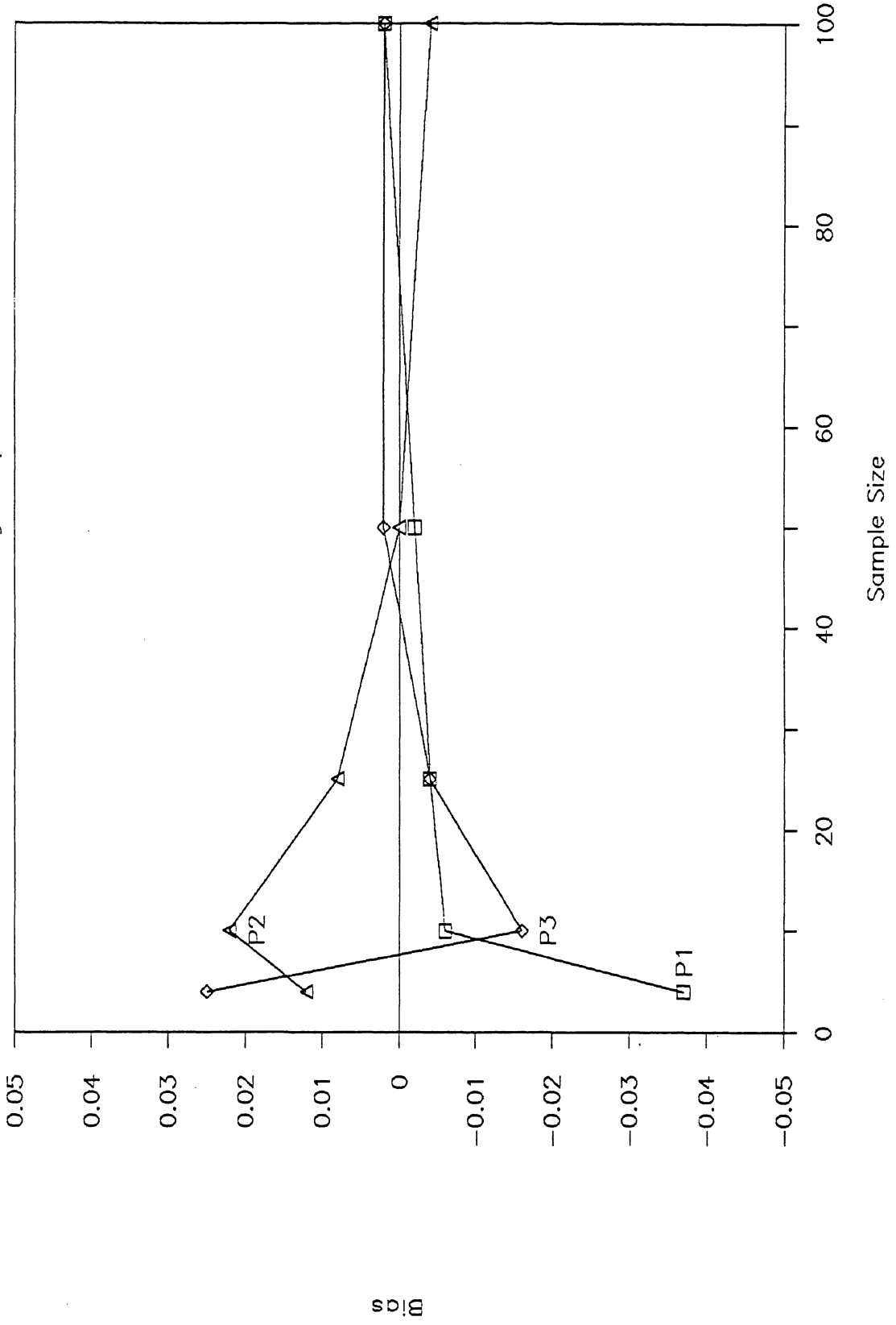


Fig 1

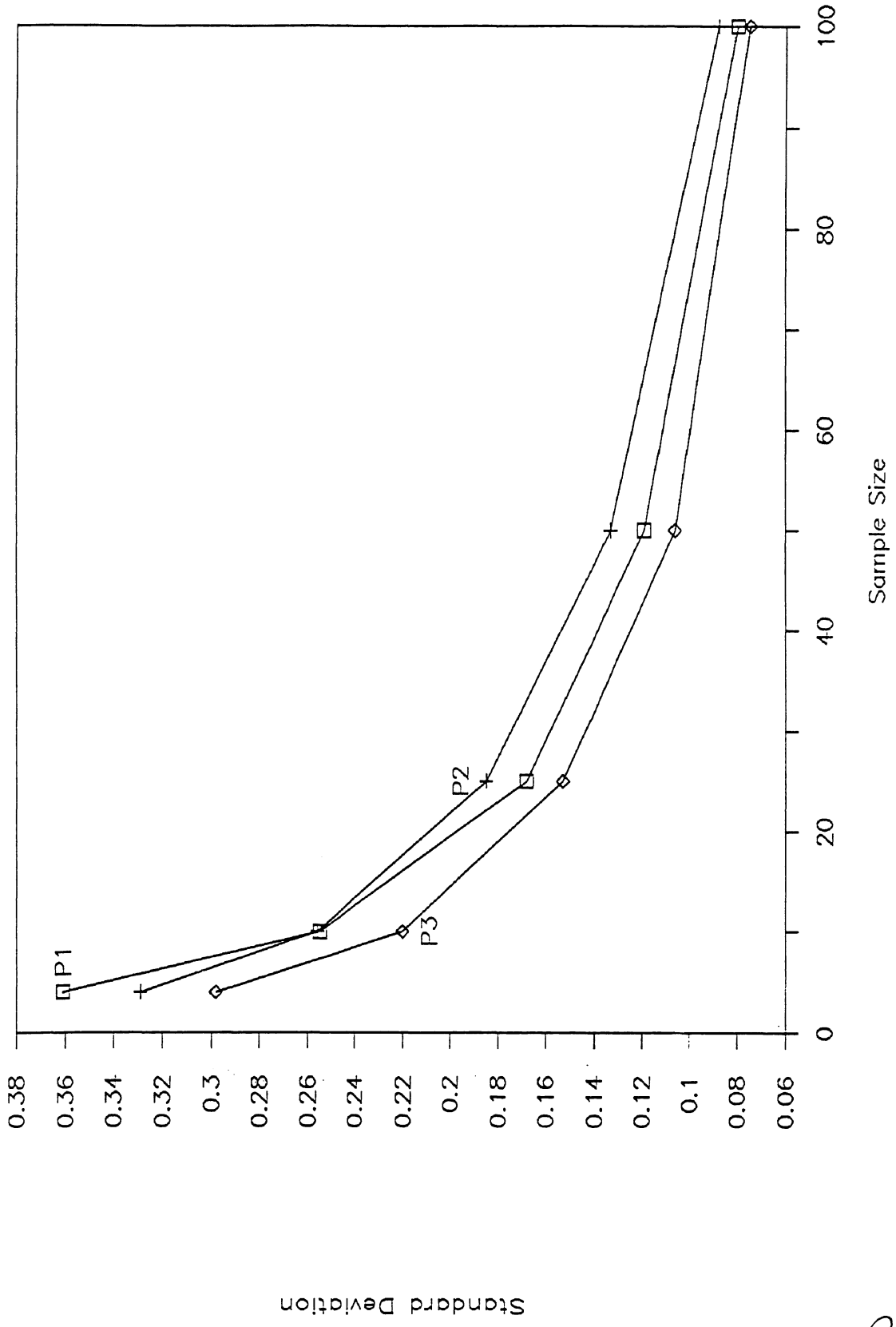


Fig 2

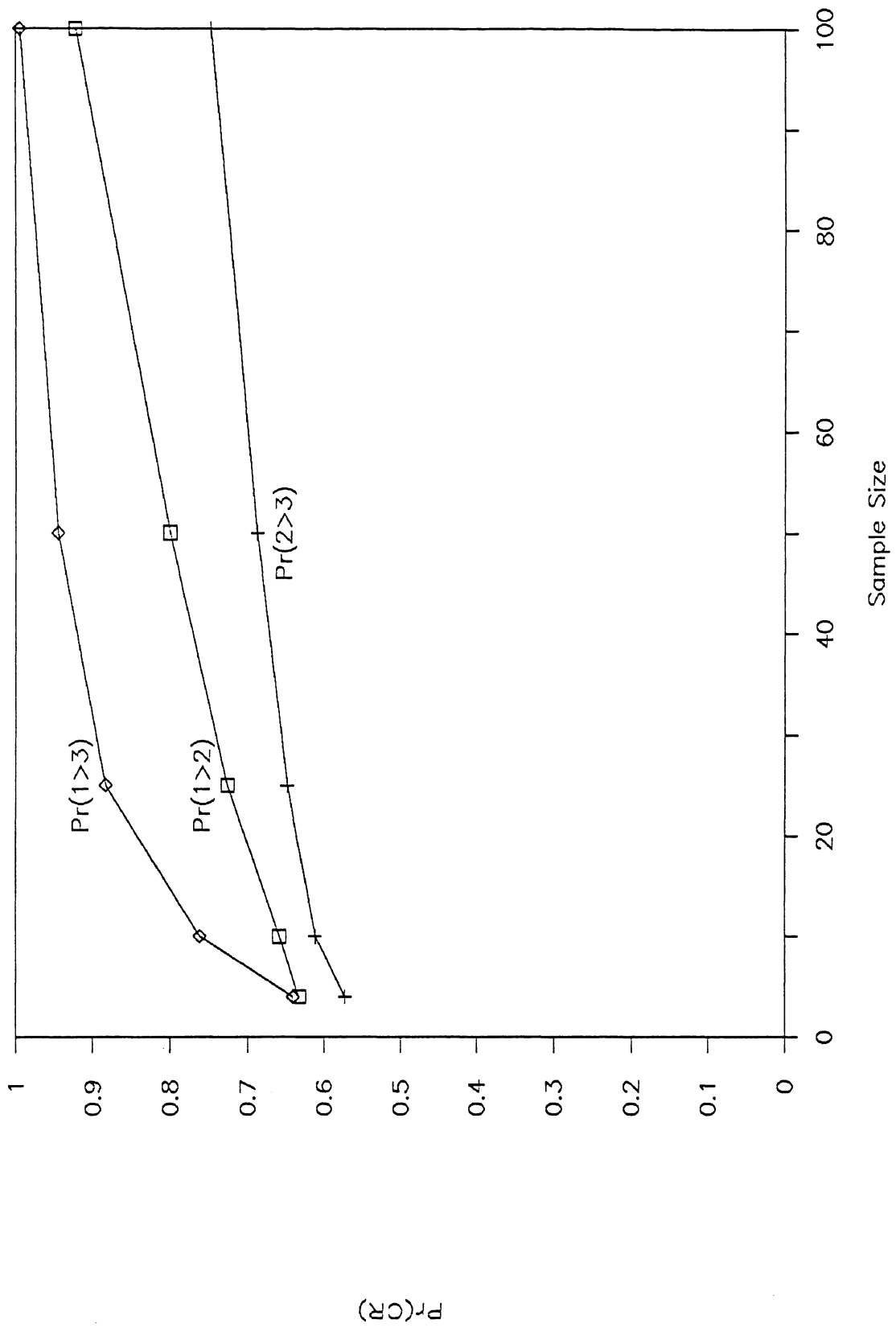


Fig 3