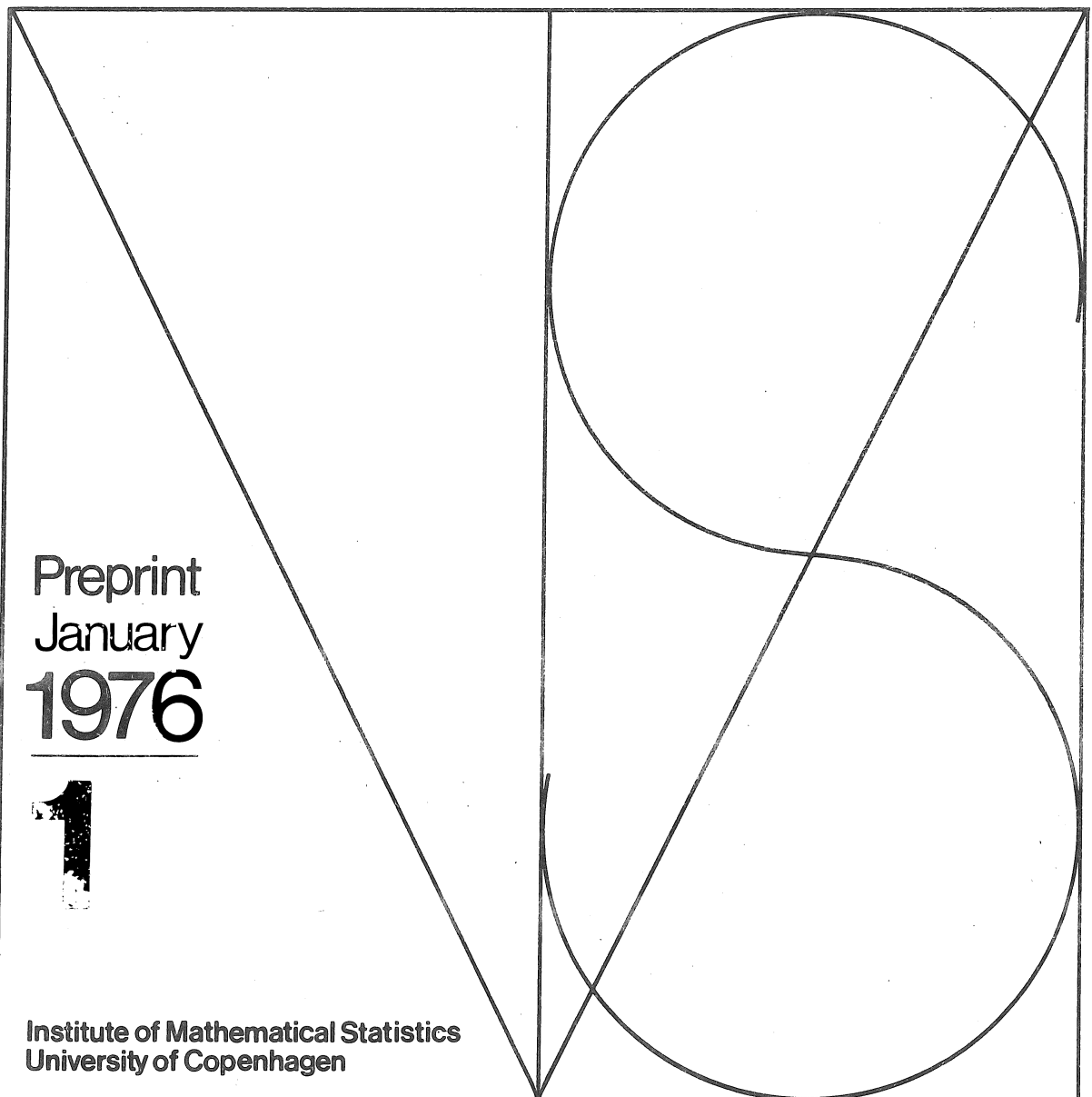


Niels Becker

Estimation for an Epidemic Model



Niels Becker^{*}

ESTIMATION FOR AN EPIDEMIC MODEL

Preprint 1976 No. 1

INSTITUTE OF MATHEMATICAL STATISTICS
UNIVERSITY OF COPENHAGEN

January 1976

* Department of Mathematical Statistic, La Trobe University, Australia.

Estimation for an epidemic model.

Niels Becker*

La Trobe University

SUMMARY

In many epidemic models the initial infection rate, suitably defined, plays a fundamental role in determining the probability of a major epidemic. An estimate for this rate is suggested on the basis of least squares and maximum likelihood estimation. The model used to arrive at the estimate is a Galton-Watson process modified by letting the offspring distribution change from generation to generation in a way as to approximate to an epidemic process. The estimates of the parameter and its associated variance are easily computed and compare well with other, computationally tedious, methods in an application to smallpox data.

Keywords: Density dependent population model; Modified Galton-Watson process; General epidemic model; Epidemic threshold theorem; Least-squares; Maximum likelihood estimation; Application to smallpox data.

1. THE PARAMETER OF INTEREST.

It is possible to formulate many mathematical models for the spread of communicable diseases in a way which clearly indicates that the initial infection rate, suitably defined, plays a fundamental role in determining the probability of a major epidemic. Recall the general epidemic model for a closed population. At time t there are $S(t)$ susceptibles, $I(t)$ infectives and $R(t)$ removals. The model is specified by the following transitions and associated probabilities for a time increment $(t, t+h)$:

<u>transition</u>	<u>probability</u>
$(S, I, R) \rightarrow (S-1, I+1, R)$	$\beta SIh + o(h)$
$\rightarrow (S, I-1, R+1)$	$\gamma Ih + o(h)$
no change	$1 - \gamma Ih - \beta SIh + o(h)$

with initial conditions $S(0) = k$, $I(0) = x_0$ and $R(0) = 0$. Whittle (1955) used this model to deduce the stochastic epidemic threshold theorem, which states essentially that for large k ,

$\text{pr}(\text{minor epidemic}) = 1 - \text{pr}(\text{major epidemic}) = \min\{1, (\gamma/\beta k)^{x_0}\}$. Here the initial infection rate $k\beta/\gamma$ determines the probability of a major outbreak. It is possible to deduce this threshold theorem from a simple birth-death process for

* The work was initiated at the Unit of Health Statistical Methodology, WHO, Geneva, and completed at the Institut for matematisk statistik, University of Copenhagen.

$\{I(t): t \geq 0\}$ by adopting $k\beta/\gamma$ as birth (death) rate and interpreting extinction (non-extinction) of the process as the event of a minor (major) epidemic. The reason is that the event of a major outbreak is essentially determined by whether or not there is a build up of infectives during the early stages of the outbreak. For large k it is reasonable to set $S(t) \approx k$ during the early stages and thereby the general epidemic model reduces to the above mentioned birth-death process. Since the simple birth-death process is an example of a branching process the following more general version of the stochastic epidemic threshold theorem suggests itself: For large susceptible populations $\text{pr}(\text{minor epidemic}) = 1 - \text{pr}(\text{major epidemic}) = q^{x_0}$, where q is the smaller root of $s = f(s)$, and f is the probability generating function of the offspring distribution for the initial infectives. Here all individuals infected by an infective are considered his offsprings. For branching processes it is known that $q < 1$ if and only if $\mu = f'(1) > 1$, and it is this result which relates the value of μ , taken as the initial infection rate in this formulation, to the probability of a major epidemic. In particular, suppose that $\mu > 1$ for a certain community and a certain disease. Suppose further that a proportion v of the community is to be vaccinated against the disease. Following the vaccination campaign the mean of the initial offspring distribution can reasonably be taken as $(1-v)\mu$. The proportion of the community that needs to be successfully vaccinated in order to prevent major outbreaks must be $\geq 1 - 1/\mu$, so that $(1-v)\mu \leq 1$. Hence the value of μ provides valuable information for deciding on disease control policies. The estimation of μ is the main concern of this paper.

The approach of Bailey and Thomas (1971) provides a method for estimating $k\beta/\gamma$ for the general epidemic model. Since $k\beta/\gamma$ may be identified with μ , their approach provides one solution to this problem. Unfortunately their method involves an impractical amount of computation. It also depends on knowing the interremoval times, which will often not be available. Finally, the dependence of their approach on the general epidemic model with its implied zero latent period and assumption of homogeneous mixing is of some concern. It is the aim here to indicate an approximate method which is computationally simple and does not depend on durations of latent or infectious periods.

The content of this paper could more generally be viewed as concerned with the estimation of the initial growth rate for a density dependent population model. Minor changes would permit application to ecological population data.

2. APPROXIMATION TO EPIDEMIC MODELS.

A branching process can approximately describe the spread of minor epidemics and the early stages of major epidemics. In order to describe the later stages or total size of major epidemics it is necessary to allow for the depletion of the susceptible population. This is done here by using a modified (embedded) Galton-Watson process to describe the development of the infective population.

The modification is to change the offspring distribution from generation to generation in a way as to allow for the depletion of the susceptible population. In the formulation it is assumed that the generation to which each infective belongs is observable. In practice this is not a realistic assumption and it is illustrated later how this difficulty can be overcome.

Let X_j denote the size of the infective population in generation j , $j=0,1,2,\dots$. The initial size X_0 is known to be x_0 . Write

$$Y_j = \sum_{i=1}^j X_i, \quad X_j = (X_1, X_2, \dots, X_j),$$

and let k denote the initial size of the susceptible population. The mean and variance of the offspring distribution in the j^{th} generation are denoted by μ_j and σ_j^2 . It is assumed throughout that $\mu_1 = \mu$, the initial infection rate, and

$$\mu_j = \mu g_j(X_{j-1}), \quad j=2,3,\dots,k,$$

where the g_j are known functions of the previous generation sizes. For homogeneously mixing populations it would be natural to take

$$g_j(X_{j-1}) = \max\{1 - Y_{j-1}/k, 0\}. \quad (1)$$

Comparisons with epidemic data has encouraged Chelsky and Angulo (1973) to consider $g_j(X_{j-1}) = (1 - Y_{j-1}/k)^2$ more appropriate. For an ecological density dependent population one might use

$$g_j(X_{j-1}) = \max\{1 - X_{j-1}/k, 0\},$$

where k is now interpreted as the saturation level of the environment.

The aim is to make as few assumptions about the offspring distribution as possible. However in order to proceed it will often be necessary to stipulate a form for σ_j^2 . It will then be assumed that

$$\sigma_j^2 = \sigma^2 h_j(X_{j-1}),$$

where the h_j are known functions of the previous generation sizes and σ^2 is an unknown parameter.

3. LEAST SQUARES ESTIMATION.

Viewing the model as an autoregressive process we write

$$X_j = \mu X_{j-1} g_j(X_{j-1}) + W_j.$$

By taking conditional expectations given X_{j-1} first, it follows that

$$E(W_j) = \text{Cov}(W_i, W_{i+j}) = 0, \quad V(W_j) = X_{j-1} \sigma_j^2.$$

To get constant variances take $\sigma_j^2 = \sigma^2 h_j(X_{j-1})$ and consider

$$X_j X_{j-1}^{-\frac{1}{2}} h_j^{-\frac{1}{2}} = \mu X_{j-1}^{\frac{1}{2}} g_j h_j^{-\frac{1}{2}} + U_j.$$

The U_j have zero mean, variance σ^2 and are uncorrelated. The least squares estimator of μ is now seen to be

$$\hat{\mu} = \frac{\Sigma(X_j g_j / h_j)}{\Sigma(X_{j-1} g_j^2 / h_j)}, \quad (2)$$

where the summation is up to n , the number of generations observed. For any stochastic process Z_1, Z_2, \dots it is known that $\Sigma[Z_j - E(Z_j | Z_{j-1})]$ is a zero mean martingale. By taking $Z_j = X_j g_j / h_j$ it follows that

$$\mu = \frac{E \Sigma(X_j g_j / h_j)}{E \Sigma(X_{j-1} g_j^2 / h_j)}$$

and hence that $\hat{\mu}$ will generally be biased. The residual sum of squares divided by $(n-1)$, namely

$$\hat{\sigma}^2 = \left[\Sigma_{j=1}^n \frac{X_j^2}{X_{j-1} h_j} - \hat{\mu} \Sigma_{j=1}^n (X_j g_j / h_j) \right] / (n-1) \quad (3)$$

suggests itself as an estimator for σ^2 . It will generally be biased. From Theorem 3 of Dion (1975) it can be deduced that the estimator $\hat{\sigma}^2$ is consistent for σ^2 when $g_j \equiv h_j \equiv 1$, for all j .

4. MAXIMUM LIKELIHOOD ESTIMATION.

Let the offspring distribution in generation j , given X_{j-1} , belong to the one-parameter exponential family with probabilities given by

$$p_j(x) = a_j(x) [c_j(\theta)]^x / A_j(\theta). \quad (4)$$

In the previous notation this implies that

$$\mu_j = \mu g_j = c_j A_j' / c_j A_j \quad \text{and} \quad \sigma_j^2 = \mu_j' c_j / c_j', \quad (5)$$

where the primes indicate derivatives with respect to θ . The log-likelihood associated with $X_{\cdot n} = x_{\cdot n}$ is then given by

$$\ell(\theta; x_{\cdot n}) = \Sigma(x_j \ln c_j - x_{j-1} \ln A_j) + \text{constant}.$$

By differentiating and using (5) one obtains

$$\frac{d\ell}{d\mu} = \Sigma(x_j - \mu g_j x_{j-1}) g_j / \sigma_j^2. \quad (6)$$

In order to obtain the maximum likelihood estimate for μ it is now necessary to specify the σ_j^2 . In particular, if $\sigma_j^2 = \sigma^2 h_j$, for all j , then the maximum likelihood

estimator for μ is the same as the least squares estimator given by (2).

The usual asymptotic properties of maximum likelihood estimators do not apply directly due to the chance of small epidemics, however, for large epidemics it seems reasonable to use the inverse of the information $\mathcal{I}(\mu) = E \left[- \frac{\partial^2 \ell(\mu; X_n)}{\partial \mu^2} \right]$ as a means of approximating the variance of the maximum likelihood estimator. From (6) it follows, after simplifying, that

$$\mathcal{I}(\mu) = E(\sum X_j g_j / \sigma_j^2) / \mu . \quad (7)$$

The results take a particularly simple form when $\sigma_j^2 = \sigma^2 g_j$, that is when the variance of the offspring distribution is proportional to its mean for each generation, namely

$$\hat{\mu} = Y_n / \sum X_{j-1} g_j , \quad \mathcal{I}(\mu) = E(Y_n) / \mu \sigma^2 . \quad (8)$$

For small epidemics the depletion of susceptibles may be ignored and so $g_j = 1$ for all j . The results of (8) then reduce to those of Becker (1974). That the results of (8) correspond to a non-trivial class of distributions even when $g_j \neq 1$ is seen by noting that $\mu_j = \mu g_j$ and $\sigma_j^2 = \sigma^2 g_j$ hold whenever

$$A_j(\theta) = [A(\theta)]^{g_j} , \quad (9)$$

where $A(\theta)$ specifies a power series distribution with probability generating function $A(z\theta)/A(\theta)$. Here μ and σ^2 are the mean and variance of the power series distribution. The family of power series distributions includes the Poisson, binomial, negative binomial and logarithmic series distributions, as well as truncated forms of these, which seems to provide adequate scope for using (8) as a basis for inference.

5. HOMOGENEOUSLY MIXING POPULATIONS.

For homogeneously mixing populations it is reasonable to take g_j as in (1). It is then possible to indicate that the model is an approximation to existing epidemic models. Although the well known Reed-Frost chain binomial model is not a branching process it is formulated in terms of "generations" and so the present notation is suitable for its description. According to the Reed-Frost model

$$\text{pr}(X_j = x_j | X_{j-1} = x_{j-1}) = \binom{k-y_{j-1}}{x_j} [1-(1-p)^{x_{j-1}}]^{x_j} (1-p)^{x_{j-1}(k-y_j)} , \quad (10)$$

where $y_j = \sum_{i=1}^j x_i$. When p is small this gives

$$E(X_j | X_{j-1} = x_{j-1}) \approx x_{j-1} (k-y_{j-1}) p \approx V(X_j | X_{j-1} = x_{j-1}) .$$

These moments are as for the above modified branching process with offspring distributions as suggested by (9) provided $\sigma^2 = \theta \frac{d\mu}{d\theta} = \mu$ holds, at least approximately.

The latter requirement is that $\mu \approx b\theta$ for some positive b . In fact $\mu = \theta$ for the Poisson offspring distribution, while $\mu \approx b\theta$ for the binomial, negative binomial and logarithmic series distribution whenever θ is small. A more careful argument for the approximation of (10) by the use of a Poisson offspring distribution is given by Ludwig (1975).

The fact that the present modified Galton-Watson process with $\mu_j = \sigma_j^2 = \rho x_{j-1}^{(k-y_{j-1})}$ is an approximation to the Reed-Frost model adds additional support to the use of

$$\hat{\mu} = \frac{ky_n}{\sum x_{j-1}^{(k-y_{j-1})}} \quad (11)$$

as a point estimate of μ , for homogeneously mixing populations.

In practice the number of cases in each generation is usually not observable. If it is known that the epidemic is terminated but only the total size y of the epidemic is observed it is useful to have bounds on $\hat{\mu}$ which involve only y , k and x_0 . With g_j as given by (1) it is known that there are at most $k+1$ generations and the denominator on the right hand side of equation (11) can then be written as

$$k(y+x_0) - \sum_{j=1}^k \sum_{i=1}^j x_i x_j. \quad (12)$$

For a given x_k consider what happens to (12) when one individual is removed from generation l and placed into generation m . The increase in (12) is found to be $x_l - x_m$. This result plus the fact that $x_i = 0 \Rightarrow x_j = 0$ for all $j > i$, implies that the expression (12) attains its maximum value, for a given total size y , when $x_j = 1, j=1, 2, \dots, y$. The corresponding minimum value is attained when $x_1 = y$. By substituting these extremes into (11) the inequalities

$$\frac{ky}{k(y+x_0) - y(y+1)/2} \leq \hat{\mu} \leq \frac{ky}{k(y+x_0) - y^2} \quad (13)$$

are obtained, for a given total size y . If $x_0 \ll y$ the inequalities simplify to

$$1/(1-y/2k) \leq \hat{\mu} \leq 1/(1-y/k).$$

When the intensity y/k of the epidemic is not too large, less than about $\frac{1}{2}$, the difference between the upper and lower bound for $\hat{\mu}$ is seen to be small. This suggests that for such epidemics there is very little information contained in the generation sizes, given the total size of the outbreak.

While the number of generations and the generation sizes are generally not observable, the duration of the epidemic is sometimes known. This information together with knowledge about the latent and infectious periods can often be used to put bounds on the number of generations. Suppose in this way it is determined that the number of generations lies strictly between l and u . The bounds in (13) can then be tightened to give

$$\frac{ky}{k(y+x_0)-uy^2/2(u-1)} \leq \hat{\mu} \leq \frac{ky}{k(y+x_0)-(y-\ell/2)^2-\ell/2} \quad (14)$$

If no information on the duration of the epidemic is available then ℓ must be taken as zero and u as $y+1$, in which case the bounds in (14) reduce to those of (13). For large outbreaks the bounds in (14) will generally not be very different from those of (13) suggesting that the duration of an epidemic does not contain much information for the estimation of μ above that contained in y .

In applications it is suggested that $\hat{\mu}$ be computed whenever it is possible to obtain the generation sizes, at least roughly. Otherwise the interval given by (13) should be used as an estimate of μ . It is then necessary to compute an associated variance. In view of (8) it seems appropriate to use $\hat{\mu}\hat{\sigma}^2/y$, where $\hat{\sigma}^2$ is given by (3). Unfortunately this approach requires the number of generations to be large and the generation sizes to be known exactly since $\hat{\sigma}^2$ is quite sensitive to variations in the generation sizes. Such data are rarely available and hence it will generally be necessary to overcome this difficulty by being more specific about the model. When the duration of the infectious period can be assumed constant and contacts between infectives and susceptibles occur randomly according to a Poisson process, then it seems appropriate to use $\hat{\mu}^2/y$ as an estimate of the variance. This suggestion is further supported by the fact that this model then provides an approximation to the Reed-Frost model. If, on the other hand, the assumptions of the general epidemic model seem more appropriate, it seems better to substitute

$$g_j(x_{j-1}) = 1-y_{j-1}/k \text{ and } \sigma_j^2 = \mu g_j(\mu g_j + 1) \quad (15)$$

into (6) and hence solve

$$\sum_{j=1}^n \frac{x_j + x_{j-1}}{k + \mu(k - y_{j-1})} = (x_0 + y_{n-1})/k \quad (16)$$

to obtain the maximum likelihood estimate $\hat{\mu}$. Then substitute (15) in (7) and hence use the reciprocal of

$$\frac{k}{\hat{\mu}^2} \sum_{j=1}^n \frac{x_j}{k + \mu(k - y_{j-1})} \quad (17)$$

as estimate of the variance of $\hat{\mu}$. The suggestion of using σ_j^2 as given by (15) is based on approximating the general epidemic model by a modified Galton-Watson process with the offspring distribution in generation j being given by (4) with

$$a_j(x) \equiv 1, \quad c_j(\theta) = \theta g_j / (1 + \theta g_j), \quad (18)$$

where g_j is as given by (1) and $\theta = k\beta/\gamma = \mu$. This is justified by assuming that during the course of an infective's infectious period the number of susceptibles

remains approximately constant and may be taken as $k-y_{j-1}$, where y_{j-1} is the number of secondary infections up to that time. The number of individuals then infected by this infective, prior to his removal, has the geometric distribution specified by (4) and (18).

Although the computations involved in the use of (16) and (17) are not too severe, it will sometimes be useful to make some further crude approximations. According to the deterministic version of the epidemic threshold theorem, see Bailey (1975), in a large epidemic the ultimate number of susceptibles will be about as much below the threshold value as it was initially above it. "On average" the infection rate will thus be about unity. Substituting $\mu_{j+1}=2$ in (15) gives

$$\mu_j = \mu(1-y_{j-1}/k), \quad \sigma_j^2 = 2\mu(1-y_{j-1}/k) \quad (19)$$

as an approximation and suggests $\hat{\mu}$ as given by (11) as an estimate of μ with $2\hat{\mu}^2/y$ as an estimate of the associated variance.

6. APPLICATION TO SMALLPOX DATA.

Bailey and Thomas (1971) quote data of D.M. Thompson and W.H. Foege on an outbreak of smallpox in a closed community in Abakaliki in south-eastern Nigeria. A total of 30 cases were observed in a population of 120 individuals at risk. The removal times, in days since the first removal, are given by 0,13, 20,22,25,25,25,26,30,35,38,40,40,42,42,47,50,51,55,55,56,57,58,60,60,61,66,66, 71,76. Using the method of maximum likelihood estimation based on the general epidemic model Bailey and Thomas, see Bailey (1975, p. 125), obtain estimates (rates per day)

$$\hat{\beta} = .00088, \quad \hat{\gamma} = .091,$$

with associated covariance matrix

$$\begin{bmatrix} 6.05 \times 10^{-8} & 5.31 \times 10^{-6} \\ 5.31 \times 10^{-6} & 9.42 \times 10^{-4} \end{bmatrix} .$$

Using the large-sample formula

$$\frac{\text{var}(\hat{\mu})}{\hat{\mu}^2} = \frac{\text{var}(\hat{\beta})}{\hat{\beta}^2} - \frac{2 \text{cov}(\hat{\beta}, \hat{\gamma})}{\hat{\beta}\hat{\gamma}} + \frac{\text{var}(\hat{\gamma})}{\hat{\gamma}^2},$$

where $\mu=k\beta/\gamma$, gives

$$\hat{\mu} = 1.15 \pm .28 .$$

When the source of infection of each case is known the generation sizes are easily deduced. Without this information it is necessary to make use of the known properties of the disease in order to obtain approximate generation sizes. It is known that the duration of the incubation period for smallpox is usually 12 days and rarely lies outside the interval 9 to 15 days. Considering clusters in the data about multiples of 12 days suggests that the generation sizes given by

generation, j	0	1	2	3	4	5	6	7
no. of cases, x _j	1	1	7	6	3	8	4	0

are reasonable. That these give the generation sizes exactly is too much to hope for, but since the estimate of μ is rather insensitive to the generation sizes this need not concern us unduly. The rough approximation to the general epidemic model obtained by using μ_j and σ_j^2 as in (19) gives

$$\hat{\mu} \pm \hat{\mu}\sqrt{2/\bar{y}} = 1.14 \pm .30 \quad (20)$$

Using instead equations (16) and (17) in order to obtain estimates for μ and the associated variance, gives

$$1.10 \pm .29 .$$

The assumptions of the general epidemic model, in particular the assumption of a zero latent period, are not very realistic for smallpox. Assuming a Poisson offspring distribution would seem to be more realistic. This leads to the same estimate of μ as in (20), namely 1.14, but the associated standard error reduces to $\hat{\mu}/\sqrt{\bar{y}} = .21$. With $x_0=1$, $k=119$ and $y=29$ the bounds in (13) reduce to $1.10 \leq \hat{\mu} \leq 1.26$. The fact that the difference between the upper and lower bounds is small illustrates that these bounds, based only on the total size y , provide a useful means of estimating μ . It also illustrates that the estimate $\hat{\mu}$ is quite insensitive to the generation sizes and that the generation sizes contain little information above that contained in y . In view of the fact that the incubation period for smallpox rarely lies outside the limits 9 to 15 days and that the duration of the outbreak was 76 days, it can be said confidently that there were at least 6 generations and at most 10 generations. Substituting $\ell=5$ and $u=11$ into (14) gives

$$1.11 \leq \hat{\mu} \leq 1.20 \quad (21)$$

Using the upper bound of (21) and, because it is the largest, the standard error of (20) it can be stated with a high degree of confidence that $\mu < 1.20 + 2 \times .30 = 1.8$. With $\mu < 1.8$ it would follow, under the assumptions of these models, that if half of the original susceptibles had been immune then a minor outbreak would have resulted with probability one.

It is interesting to note that a reasonable estimate of μ is obtained by using the deterministic general epidemic model. In the present notation the deterministic model, see Bailey (1975, equation 6.14), gives

$$\mu = - \frac{x_0+k}{x_0+y} \ln(1-y/k) \quad (22)$$

Substituting $x_0=1$, $k=119$ and $y=29$ leads to the estimate 1.12 of μ . This deterministic approach of course does not provide any indication as to the accuracy of the estimate. In an announcement Startsev (1970) suggests

$$\frac{x_0 + y}{k [N(k) - N(k - y - 1)]} \quad (23)$$

where $N(k) = \sum_{j=1}^k 1/j$, as an estimate of $1/\mu$ for the stochastic general epidemic model. Since $N(n) - \ln n \approx \text{constant}$ for large n , it follows that (23) gives the same estimate as (22) whenever $x_0 \ll k$ and $y \ll k$. These conditions are also used by Startsev to justify (23) as an estimate.

By way of summary it is noted that the estimates provided by the methods of the present paper are in close agreement with those of Bailey and Thomas. The main advantage of the present approach lies in the simplicity of the computations involved and the fact that some of the results involve only the total size of the outbreak.

REFERENCES

- BAILEY, N.T.J. (1975). The mathematical theory of infectious diseases. 2nd edition. Griffin. London.
- BAILEY, N.T.J. and THOMAS, A.S. (1971). The estimation of parameters from population data on the general stochastic epidemic. Theor. Pop. Biol., 2, 253-70.
- BECKER, N.G. (1974). On parametric estimation for mortal branching processes. Biometrika, 61, 393-9.
- CHELSEY, M. and ANGULO, J.J. (1973). Two models for estimation of some parameters of disease spread. Math. Biosci., 18, 119-31.
- DION, J.-P. (1975). Estimation of the variance of a branching process. Ann. Statist., 3, 1183-7.
- LUDWIG, D. (1975). Qualitative behaviour of stochastic epidemics. Math. Biosci., 23, 47-73.
- STARTSEV, A.N. (1970). On estimation of the regulating parameter in the stochastic model for epidemics. Kratkie Nauchnye Soobstcheniya. No. 6. Izvestija Akademia Nauk Uzbekskoj SSR 1970.
- WHITTLE, P. (1955). The outcome of a stochastic epidemic - a note on Bailey's paper. Biometrika, 42, 116-22