# Non-parametric estimation of the case fatality ratio with competing risks data: An application to Severe Acute Respiratory Syndrome (SARS)

Nicholas P. Jewell[1, *, †], Xiudong Lei[1, ‡], Azra C. Ghani[2, §], Christl A. Donnelly[2, ¶], Gabriel M. Leung[3, ‖], Lai-Ming Ho[3, **], Benjamin J. Cowling[3, ††] and Anthony J. Hedley[3, ‡‡]

[1]*Division of Biostatistics and Department of Statistics, University of California, Berkeley, U.S.A.*
[2]*Department of Infectious Disease Epidemiology, Imperial College, London, U.K.*
[3]*Department of Community Medicine, University of Hong Kong, Hong Kong*

## SUMMARY

For diseases with some level of associated mortality, the case fatality ratio measures the proportion of diseased individuals who die from the disease. In principle, it is straightforward to estimate this quantity from individual follow-up data that provides times from onset to death or recovery. In particular, in a competing risks context, the case fatality ratio is defined by the limiting value of the sub-distribution function, $F_1(t) = \Pr(T \leqslant t \text{ and } J = 1)$, associated with death, as $t \to \infty$, where $T$ denotes the time from onset to death ($J = 1$) or recovery ($J = 2$). When censoring is present, however, estimation of $F_1(\infty)$ is complicated by the possibility of little information regarding the right tail of $F_1$, requiring use of estimators of $F_1(t^*)$ or $F_1(t^*)/(F_1(t^*) + F_2(t^*))$ where $t^*$ is large, with $F_2(t) = \Pr(T \leqslant t \text{ and } J = 2)$ being the analogous sub-distribution function associated with recovery. With right censored data, the variability of such estimators increases as $t^*$ increases, suggesting the possibility of using estimators at lower values of $t^*$ where bias may be increased but overall mean squared error be smaller. These issues are investigated here for non-parametric estimators of $F_1$ and $F_2$. The ideas are illustrated on case fatality

*Correspondence to: Nicholas P. Jewell, Division of Biostatistics and Department of Statistics, University of California, Berkeley, U.S.A.
†E-mail: jewell@stat.berkeley.edu
‡E-mail: xiudonglei@hotmail.com
§E-mail: azra.ghani@lshtm.ac.uk
¶E-mail: c.donnelly@imperial.ac.uk
‖E-mail: gmleung@hku.hk
**E-mail: lmho@hkucc.hku.hk
††E-mail: bcowling@hkucc.hku.hk
‡‡E-mail: hrmrajh@hkucc.hku.hk

data for individuals infected with Severe Acute Respiratory Syndrome (SARS) in Hong Kong in 2003. Copyright © 2006 John Wiley & Sons, Ltd.

## 1. INTRODUCTION

The case fatality ratio (CFR) of a disease measures the proportion of afflicted individuals who die from the disease. For infectious or other acute diseases, the CFR is an important measure of virulence, and is often used to assess the effect of cofactors and intervention strategies. In the early stages of an epidemic, estimation of the CFR is complicated by the fact that the disease may not have run its course in many affected individuals at the time of analysis; that is, the times to death or recovery are right censored. The implications of this are best understood in the context of a competing risks model.

Suppose the outcome of a disease is either death or recovery, and let the random variables $T$ and $J$ measure the time from initiation (e.g. infection) until the final outcome, and result of outcome (say $J = 1$ corresponds to death, and $J = 2$ to recovery), respectively. The two sub-distribution functions of primary interest are

$$F_j(t) = \mathrm{pr}(T \leqslant t, J = j), \quad j = 1, 2$$

with the overall survival function given by

$$S(t) = \mathrm{pr}(T > t) = 1 - F_1(t) - F_2(t)$$

measuring the probability of neither recovering nor dying by time $t$. (Note that the standard survival analysis terminology is unfortunate here since 'survival' up to time $t$ in this context means that no event has occurred by time $t$, and not recovery.) The CFR is then simply defined by $\lim_{t \to \infty} F_1(t) = \mathrm{pr}(J = 1)$, the proportion of diseased individuals who eventually die; for simplicity we refer to the CFR by $F_1(\infty)$.

Suppose that, for each individual, information on survival status is subject to independent right censoring at time $C$. Thus, the observed data can be represented as $Y = (\tilde{T}, \Delta, \Delta J)$, where $\tilde{T} = \min(T, C)$ and $\Delta = 1$ if $T < C$ and is 0 otherwise. Note that observation of the product term $\Delta J$ simply means that the cause of failure—either death or recovery—is known whenever an uncensored event is observed, but not otherwise. We assume that the censoring random variable $C$ follows an unknown distribution function $G$. In Section 2 we review non-parametric procedures for estimation of $F_1(t)$ and $F_2(t)$, leading to estimators $\hat{F}_1(t)$ and $\hat{F}_2(t)$. The presence of right censoring complicates non-parametric estimation of $F_1(\infty)$ as information on the value of $F_1$ may run out before the sub-distribution function is close to its asymptotic limit. Given this obstacle, two estimation strategies have been suggested for $F_1(\infty)$. First, we can estimate $F_1(\infty)$ by

$$\mathrm{CFR}_a = \hat{F}_1(t^*)$$

for a suitably chosen large value $t^*$ (perhaps, the largest observed death time). An alternative estimator is given by

$$\text{CFR}_b = \frac{\hat{F}_1(t^*)}{\hat{F}_1(t^*) + \hat{F}_2(t^*)}$$

for a similarly large $t^*$ (perhaps, the largest observed outcome time, whether death or recovery). The first of these estimators ignores the fact that individuals may die after time $t^*$, and the second assumes that the proportion of individuals, still afflicted at time $t^*$, who ultimately die is the same as observed for those whose outcome occurs prior to $t^*$. Gaynor *et al.* [1] briefly discuss these estimators. In Section 2, we introduce two familiar non-parametric strategies for estimation of $F_1$ and $F_2$, and discuss various methods for estimation of the (asymptotic) variances of both $\text{CFR}_a$ and $\text{CFR}_b$.

These two strategies for estimation of the CFR are analogies to two estimators based on group rather than individual data; that is, where only the total number of deaths and recoveries are known by a fixed follow-up or censoring time $C$. The naive estimate of the CFR as the fraction of all afflicted individuals who die by time $C$ is clearly subject to bias, particularly if the fixed time $C$ is small relative to the upper bound of the support of $T$. The second approach uses the fraction of deaths amongst those whose outcome is known at time $C$. For grouped SARS data, the second approach was recommended by several authors [2, 3].

For either $\text{CFR}_a$ or $\text{CFR}_b$, choosing $t^*$ large enough is important in reducing bias; however, in the presence of substantial right censoring, non-parametric estimators of $\hat{F}_1$ and $\hat{F}_2$ suffer from increased variability for large $t^*$. These two observations suggest that the value $t^*$ might be chosen somewhat smaller, trading an increase in bias for decreased variance in order to reduce mean squared error. In Section 2.2 we suggest a data adaptive procedure for the selection of $t^*$. In Section 3 we use some limited simulations to assess the performance of estimators of $\text{CFR}_a$ and $\text{CFR}_b$, and also compare different estimators of their sampling variability, namely (i) the (asymptotic) estimated variance based on influence curve calculations, (ii) a Greenwood-type estimator of the asymptotic variance, (iii) a simple bootstrap variance estimate, and (iv) an asymptotic approximation suggested by Cox [2, 4]. Finally, in Section 4, the ideas are illustrated on data from SARS patients in Hong Kong in Spring 2003 [2, 4].

## 2. NON-PARAMETRIC ESTIMATORS OF $F_1$ AND $F_2$ AND THEIR ASYMPTOTIC VARIANCES

Two basic strategies are available for non-parametric estimation of $F_1$ and $F_2$. To describe and compare these we first require some notation. Let $t_1 < \cdots < t_k$ denote the distinct observed event times for outcomes of either type, with $d_{ij}$ representing the number of outcomes of type $j$ that occur at time $t_i$, and $n_i$ the number of subjects at risk at time $t_i$. The first estimator is non-parametric maximum likelihood, yielding

$$\hat{F}_1(t)_{\text{ML}} = \sum_{t_i \leqslant t} \frac{d_{i1}}{n_i} \hat{S}(t_i^-) \tag{1}$$

with an analogous expression for $\hat{F}_2(t)_{\text{ML}}$ (Chapter 8.2 [5]), where $\hat{S}$ is the Kaplan–Meier estimator of the overall survival function.

For complete data $(T, J)$, the non-parametric maximum likelihood estimator of $F_1$ is the empirical sub-distribution function given by $n^{-1} \sum_{i=1}^{n} I(T_i \leqslant t, J = 1)$, where $I(E)$ is the

indicator function for the event $E$ and $n$ is the total sample size. With censored data this suggests using this estimator for the observed complete (that is, uncensored) observations, weighted by the probability of not being censored. This *inverse probability of censoring weighted estimator* is therefore given by

$$\hat{F}_1(t)_{\text{IPCW}} = \frac{1}{n} \sum_{i=1}^{n} \frac{I(t_i \leqslant t, J = 1)\delta_i}{1 - \hat{G}(t_i)}$$

where $\hat{G}$ is an estimator of the censoring distribution function $G$, say, the non-parametric maximum likelihood (Kaplan–Meier) estimator. Again, an analogous definition yields $\hat{F}_2(t)_{\text{IPCW}}$. Straight-forward algebra establishes that, in fact, $\hat{F}_1(t)_{\text{ML}} = \hat{F}_1(t)_{\text{IPCW}}$, with a similar equivalence for the two estimators of $F_2$. From this point, we therefore refer to either estimator as $\hat{F}_1(t)$ and, similarly, $\hat{F}_2(t)$. Inversely weighting uncensored observations by the (estimated) probability of censoring has a long history in survival analysis; it was employed for the censored linear regression problem by Koul *et al.* [6], and in a much more general setting by Robins and Rotnitzky [7].

The estimator (1) and its analogue for $F_2$ are special cases of the Aalen–Johansen estimator [8], previously developed for competing risks data by Aalen [9], and have been discussed extensively in the competing risks literature. For example, Gooley *et al.* [10] develop these estimators using the redistribute to the right algorithm.

## 2.1. Asymptotic variance estimation for CFR$_a$ and CFR$_b$

The influence curve is a useful concept allowing some estimators to be approximately expressed in terms of an average of component estimators each of which depend on a single observation. Among its many uses, this can allow straightforward determination of the large sample properties of the estimator, particular the asymptotic variance. An introduction to influence curves is given in Reid [11].

Based on the observed data $Y = (\tilde{T}, \Delta, \Delta J)$, it can be shown that $\hat{F}_1(t)$ is an asymptotically linear estimator of $F_1(t)$ with influence curve $\text{IC}_1(Y; t)$ meaning that $\hat{F}_1(t) - F_1(t)$ can be approximated by an empirical mean of $\text{IC}_1(Y; t)$

$$\hat{F}_1(t) - F_1(t) = \frac{1}{n} \sum_{i=1}^{n} \text{IC}_1(Y; t) + o_P(1/\sqrt{n})$$

In Appendix A, we provide the formula for the influence curve $\text{IC}_1(Y; t)$ and establish that $\sqrt{n}(\hat{F}_1(t) - F_1(t))$ converges in distribution to a normal distribution with mean zero and variance $\sigma_1^2 = E\{\text{IC}_1(Y; t)^2\}$. In principal, this asymptotic variance can be estimated consistently with $\hat{\sigma}_1^2 = \frac{1}{n} \sum_{i=1}^{n} \hat{\text{IC}}_1(Y_i; t)^2$, where $(T_i, J_i)$, and non-parametric maximum likelihood estimators, $\hat{F}_1$ and $\hat{G}$, are substituted in (A1) for $(T, J)$, $F_1$ and $G$, respectively, to obtain $\hat{\text{IC}}_1(Y_i; t)$; asymptotic confidence intervals can then be constructed in the standard fashion, perhaps with use of a transformation to improve the approximation in moderate samples. The influence curve, $\text{IC}_2(Y; t)$, and the asymptotic variance, for $\hat{F}_2(t)$, are determined and estimated in an analogous way.

It is well-known however that plug-in estimators of the variance of an influence curve may not perform well in finite samples. In the simple case of competing risks, there is a more effective estimator of $\sigma_1^2$ given by a generalization of Greenwood's formula to the competing risk setting [8] which is described briefly in Appendix B (see also Example IV.4.1, pp. 298–304 [12]).

A somewhat similar estimator, that is nevertheless very similar numerically to the Greenwood-type method, was originally developed by Dinse and Larson [13] in the context of a semi-Markov model—see also Gaynor *et al.* [1]. Pepe [14] and Lin [15] discuss alternative variance estimators that are more complex computationally. Choudhury [16] briefly describes and compares these various estimators of the asymptotic variance via limited simulations, finding that a Greenwood estimator, such as (B2) works well even with relatively small sample sizes, and that coverage probabilities for associated confidence intervals are also satisfactory so long as $n > 200$. It again may be valuable to use a transformation (e.g. $\log - \log\{\hat{F}_1(t)\}$) to improve the asymptotic approximation underlying such confidence intervals, particularly for small sample sizes.

We now turn briefly to these same issues associated with estimation of CFR$_b$. As described in Appendix A, the influence curve, $\text{IC}(Y; t)$, for CFR$_b$ is simply obtained from those for $\hat{F}_1(t)$ and $\hat{F}_2(t)$. As before, we can estimate the asymptotic variance $\sigma^2$ of $\hat{F}_1(t)/(\hat{F}_1(t) + \hat{F}_2(t))$ by $\hat{\sigma}^2 = \frac{1}{n}\sum_{i=1}^{n}\hat{\text{IC}}(Y_i; t)^2$ where an estimate of $\text{IC}(Y_i; t)$ is obtained by plugging in empirical estimates of the various components of (A2) given in Appendix A. Again, a more reliable estimator of the asymptotic variance of CFR$_b$ can be based on the Greenwood formulae (B1) and (B2) together with application of the delta method to yield

$$\widehat{\text{var}}\left(\frac{\hat{F}_1(t)}{\hat{F}_1(t)+\hat{F}_2(t)}\right) = \frac{[\hat{F}_2(t)]^2\,\text{var}(\hat{F}_1(t))+[\hat{F}_1(t)]^2\,\text{var}(\hat{F}_2(t))-2\hat{F}_1(t)\hat{F}_2(t)\widehat{\text{cov}}(\hat{F}_1(t), \hat{F}_2(t))}{(\hat{F}_1(t) + \hat{F}_2(t))^4}$$

(2)

As for CFR$_a$, the use of a transformation may improve the coverage probability of a confidence interval based on (2) for small sample sizes.

Variance estimators such as Greenwood's formula, approximate the variance of $\hat{F}_1(t)$ and the ratio, $\hat{F}_1(t)/(\hat{F}_1(t) + \hat{F}_2(t))$, at a fixed time $t$. However, the estimate CFR$_a$ is based on $\hat{F}_1(t^*)$ at a perhaps randomly selected time $t^*$, as when using the largest observed death time. Similarly, CFR$_b$ is evaluated at a randomly selected $t^*$. The asymptotic variance estimators do not account for variability associated with determination of $t^*$. An alternative estimator of the variance that does allow for the random choice of $t^*$ is based on the bootstrap technique. For specific implementation, bootstrap samples of $Y = (\tilde{T}, \Delta, \Delta J)$ of size $n$ are randomly selected, with replacement, from the original data. The choice of the number of bootstrap samples is delicate and is discussed in detail in Shao and Tu [17]. As a rule of thumb they suggest 50–200 bootstrap samples for moment estimators like the variance used here. In the simulations of Section 3 and the example in Section 4, we therefore use 200 bootstrap samples. If computational considerations are not of concern, the number of bootstrap samples can be increased until stability in the variance estimates is achieved.

The influence curves of $\hat{F}_1(t)$ and $\hat{F}_2(t)$, and thus their variance, involve the term $\bar{G}(\cdot)$ in the denominators as in (A1). We can thus expect the asymptotic variances of $\hat{F}_1(t)$ and $\hat{F}_2(t)$ to be large in the 'tail' where $\bar{G}(t)$ is close to zero. Unfortunately, minimizing bias suggests choosing $t^*$ as large as possible when calculating CFR$_a$ or CFR$_b$. However, the possible increased variability motivates using a smaller value of $t^*$. This is explored in Section 2.3 and in the simulations in Section 3. Finally, we note that, in the absence of censoring, the influence curve method and Greenwood formulae such as (A2) reduce to standard variance estimators based on the simple multinomial distribution for outcome counts by a fixed time.

Finally, by ignoring certain correlation terms, Cox suggested a clever approximation to the asymptotic variance of CFR$_b$ (see Ghani *et al.* [4]) that the latter authors used as an alternative to

either the influence curve or Greenwood-type variance estimators. The details of Cox's approximation can also be found in Appendix C.

## 2.2. Choosing the value of $t^*$

To minimize bias, it is natural to use the largest observed death and/or recovery time for $t$ in constructing the estimators $CFR_a$ and $CFR_b$; to emphasize the dependence on the choice of $t$ in these estimators, we here denote them, respectively, by $CFR_a(t_{max})$ and $CFR_b(t_{max})$. On the other hand, at these observed values of $t$, the associated standard errors may be large, as the probability of censoring is high. It is plausible that the mean squared error of these estimators might be improved by selecting a smaller value of $t$ in $CFR_a(t)$ and $CFR_b(t)$ in order to decrease variability at the cost of admitting some extra bias. Further, it is possible to consider a data-adaptive choice of an 'optimal' $t$, selected to minimize mean squared error in this way.

Suppose, for example, that $t_{max}$ is the largest observed death time. For $t < t_{max}$, the mean squared error of $\hat{F}_1(t)$ as an estimate of CFR can, in turn, be estimated by

$$\{\hat{F}_1(t) - \hat{F}_1(t_{max})\}^2 + \{\hat{\sigma}_1^2(t)\} \tag{3}$$

where $\hat{\sigma}_1^2(t)$ estimates the variance of the estimator of Section 2.1 at the value $t$. The 'optimal' time to estimate CFR, $t_{opt}$, is then selected to be that value of $t$ that empirically minimizes (3).

A similar approach can be used to choose an optimal $t$ when using $CFR_b$; that is we seek the $t$ that minimizes the estimated mean square error of $CFR_b$

$$\left\{ \frac{\hat{F}_1(t)}{\hat{F}_1(t) + \hat{F}_2(t)} - \frac{\hat{F}_1(t_{max})}{\hat{F}_1(t_{max}) + \hat{F}_2(t_{max})} \right\}^2 + \{\hat{\sigma}^2(t)\} \tag{4}$$

where $\hat{\sigma}^2(t)$ is an appropriate estimator of the variance of $CFR_b$ as discussed in Section 2.1.

It is likely that the optimal time values, for estimation of $CFR_a$ and $CFR_b$, differ. We emphasize that use of variance estimates, such as $\hat{\sigma}_1^2(t_{opt})$, ignores variability associated with the random choice of $t_{opt}$. The simulations of Section 3 provide insight into whether use of $t_{opt}$, in place of $t_{max}$, is likely to be of value.

## 3. SIMULATIONS

A limited simulation study was employed to compare the performance of $CFR_a$ and $CFR_b$, and to assess the value in estimating the case fatality ratio at $t_{opt}$ as compared to $t_{max}$. We consider three different scenarios, characterized by differing censoring patterns. In the motivating example of an epidemic, the censoring distribution is generated by the arrival pattern in time of newly infected cases in relation to the chronological time of data analysis. The three scenarios thus, in part, reflect analysis of outcome data at differing times in the course of an epidemic. We specify the joint distribution of $(T, J)$ through the CFR and the two conditional distributions of $T$, given $J = 1$ and given $J = 2$.

The simulation parameters were motivated by the parametric approximation of SARS data used by Donnelly et al. [2]. For each scenario, the true CFR is set to be 0.2, with conditional outcome distributions both Gamma, with means 35 and 25 for death and recovery, respectively, with both having variance 200. In scenario I, the censoring distribution is Uniform on [0, 100], yielding an

overall probability of censoring of about 30 per cent. In scenario II, the censoring distribution is a mixture of a Uniform distribution on (0, 50) (with probability 0.2) and, beyond 50, an Exponential distribution (with origin at 50) with a rate parameter of 0.2 (with probability 0.8)—this yields an overall probability of censoring of about 19 per cent. Finally the censoring distribution for scenario III is similar to that for scenario II, except that the change point between the Uniform and Exponential distributions occurs at 30 rather than 50, and beyond 30 the Exponential distribution has rate parameter 0.1, these changes yielding a higher overall censoring probability of about 36 per cent. For the three scenarios, Figures 1 and 2 display the survival distribution corresponding to censoring, $\bar{G}(t)$, together with $F_1(t)$ and $F_1(t)/(F_1(t) + F_2(t))$ for scenarios I and III, respectively. With regard to data arising from an epidemic, scenario I is intended to reflect an analysis at a mature state of the epidemic with a substantial amount of complete information on death or recovery; scenario III, on the other hand, mimics an analysis much earlier in an epidemic so that many observations are censored before information is available on death or recovery. Scenario II is intermediate between these two situations.

For each scenario 1000 data sets of two extremes of sample size were simulated, one with $n = 100$, and the other with $n = 1500$, the latter case roughly corresponding with the Hong Kong SARS data considered in Section 4. For each simulation, the estimators, $\text{CFR}_a$ and $\text{CFR}_b$ at $t_{\max}$, and $\text{CFR}_b$ at the data-driven optimal $t_{\text{opt}}$ (as described in Section 3) and their estimated variances (based on the appropriate influence curve, Greenwood formula, Cox approximation or bootstrap for $\text{CFR}_a$ and $\text{CFR}_b$; only the Greenwood formula and Cox approximation for $\text{CFR}_b$ at $t_{\text{opt}}$ for computational reasons) were computed. Either the Greenwood formula or the Cox approximation was used in evaluating estimated mean squared error in choosing the value $t_{\text{opt}}$ for each data set. For the bootstrap variance estimator, 200 bootstrap samples with replacement were generated; since $t_{\max}$ will vary across these pseudo-samples, this method accounts for this form of variation unlike the other methods. For the two sample sizes, $n = 100$ and 1500, the simulation mean and variance of three estimators are reported in Tables I and II, respectively, along with the simulation
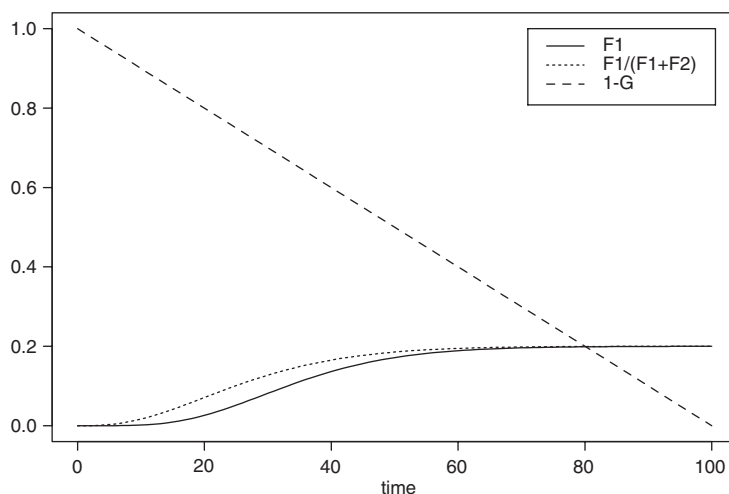


Figure 1. Description of the sub-distribution function $F_1$ and ratio $F_1/(F_1 + F_2)$ and censoring distribution $G$ for simulation scenario I.
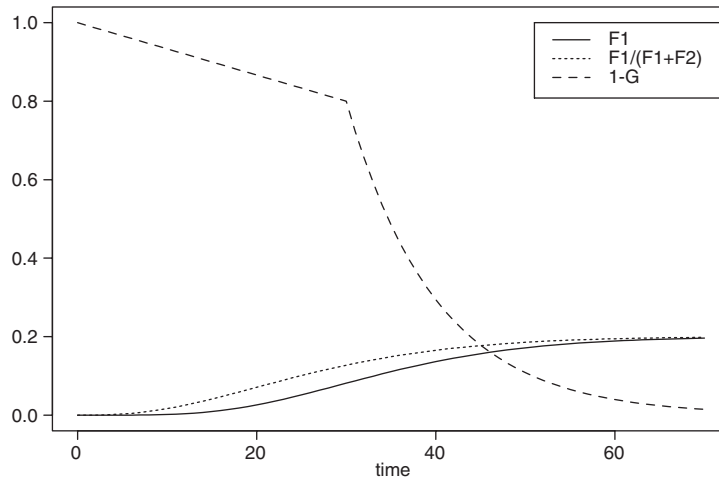
Figure 2. Description of the sub-distribution function $F_1$ and ratio $F_1/(F_1 + F_2)$ and censoring distribution $G$ for simulation scenario III.

Table I. Comparison of estimators of the case fatality ratio in three simulation scenarios ($n = 100$).

| Scen. | Estimator | Simul. mean | Simul. variance ($\times 10^{-2}$) | Simul. MSE ($\times 10^{-2}$) | Simulation mean of variance estimators ($\times 10^{-2}$) (95% CI Coverage) | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | IC | bootstrap | Cox | Greenwood |
| I | $\mathrm{CFR}_a(t_{\max})$ | 0.1855 | 0.2828 | 0.3038 | 0.2821 (90.6) | 0.2733 (90.6) | 0.4435 (93.0) | 0.2418 (89.4) |
| | $\mathrm{CFR}_b(t_{\max})$ | 0.1912 | 0.2841 | 0.2918 | 0.2410 (88.4) | 0.2739 (91.4) | 0.2871 (89.4) | 0.3452 (94.8) |
| | $\mathrm{CFR}_b(t_{\mathrm{opt}})$ | 0.1764 | 0.2611 | 0.3168 | | | 0.1040 (68.8) | |
| | | 0.1756 | 0.2937 | 0.3532 | | | | 0.2498 (83.0) |
| II | $\mathrm{CFR}_a(t_{\max})$ | 0.1550 | 0.2612 | 0.4637 | 0.4124 (69.2) | 0.2061 (68.0) | 0.3558 (75.6) | 0.1756 (65.8) |
| | $\mathrm{CFR}_b(t_{\max})$ | 0.1644 | 0.2555 | 0.3822 | 0.3147 (78.8) | 0.2138 (78.6) | 0.2549 (77.6) | 0.2832 (88.0) |
| | $\mathrm{CFR}_b(t_{\mathrm{opt}})$ | 0.1480 | 0.2033 | 0.4737 | | | 0.0693 (43.2) | |
| | | 0.1462 | 0.2236 | 0.5130 | | | | 0.1902 (68.6) |
| III | $\mathrm{CFR}_a(t_{\max})$ | 0.1486 | 0.7459 | 1.0101 | 1.1632 (68.4) | 0.4613 (64.2) | 0.7380 (63.4) | 0.3039 (58.2) |
| | $\mathrm{CFR}_b(t_{\max})$ | 0.1583 | 0.7270 | 0.9009 | 0.6031 (70.4) | 0.4603 (69.2) | 0.4764 (66.4) | 0.3887 (70.6) |
| | $\mathrm{CFR}_b(t_{\mathrm{opt}})$ | 0.1353 | 0.5348 | 0.9534 | | | 0.0780 (32.8) | |
| | | 0.1328 | 0.5628 | 1.0144 | | | | 0.1688 (45.6) |

Table II. Comparison of estimators of the case fatality ratio in three simulation scenarios ($n = 1500$).

| Scen. | Estimator | Simul. mean | Simul. variance ($\times 10^{-3}$) | Simul. MSE ($\times 10^{-3}$) | Simulation mean of variance estimators ($\times 10^{-3}$) (95% CI Coverage) | | | |
|-------|-----------|-------------|-------------------------------------|--------------------------------|------|-----------|------|-----------|
| | | | | | IC | bootstrap | Cox | Greenwood |
| I | $\mathrm{CFR}_a(t_{\max})$ | 0.1955 | 0.198 | 0.218 | 10.188 (100) | 0.194 (94.2) | 0.233 (95.4) | 0.168 (91.4) |
| | $\mathrm{CFR}_b(t_{\max})$ | 0.1966 | 0.194 | 0.206 | 1.625 (100) | 0.191 (94.6) | 0.153 (91.2) | 0.239 (98.0) |
| | $\mathrm{CFR}_b(t_{\mathrm{opt}})$ | 0.1948 | 0.190 | 0.217 | | | 0.127 (86.0) | |
| | | 0.1953 | 0.190 | 0.212 | | | | 0.230 (96.2) |
| II | $\mathrm{CFR}_a(t_{\max})$ | 0.1772 | 0.577 | 1.097 | 26.347 (100) | 0.411 (64.0) | 0.254 (56.6) | 0.192 (48.8) |
| | $\mathrm{CFR}_b(t_{\max})$ | 0.1822 | 0.484 | 0.801 | 12.964 (97.0) | 0.355 (71.0) | 0.180 (58.2) | 0.245 (66.0) |
| | $\mathrm{CFR}_b(t_{\mathrm{opt}})$ | 0.1775 | 0.374 | 0.880 | | | 0.118 (43.6) | |
| | | 0.1776 | 0.372 | 0.874 | | | | 0.202 (57.6) |
| III | $\mathrm{CFR}_a(t_{\max})$ | 0.1815 | 0.997 | 1.339 | 117.507 (100) | 0.804 (78.8) | 0.340 (63.8) | 0.327 (63.2) |
| | $\mathrm{CFR}_b(t_{\max})$ | 0.1856 | 0.879 | 1.086 | 3.912 (100) | 0.728 (83.2) | 0.247 (63.0) | 0.367 (72.2) |
| | $\mathrm{CFR}_b(t_{\mathrm{opt}})$ | 0.1817 | 0.795 | 1.130 | | | 0.147 (50.6) | |
| | | 0.1815 | 0.802 | 1.144 | | | | 0.207 (62.4) |

mean of corresponding variance estimators and the empirical coverage of nominal 95 per cent confidence intervals based on each variance estimator.

Generally, $\mathrm{CFR}_a$ exhibits a little more bias and variability than $\mathrm{CFR}_b$, as can be expected from the shapes of $F_1$ and $F_1/(F_1 + F_2)$ as illustrated in Figures 1 and 2. For all estimators, the bias is worst in scenarios II and III where it is considerable due to heavier censoring in the tail of $F_1$ and $F_2$. The estimator, $\mathrm{CFR}_b$, evaluated at $t_{\mathrm{opt}}$, is more biased as expected, but also displays the anticipated gains in precision. Figure 3 shows the behaviour of bias squared, variance, and mean squared error for a typical data set with $n = 100$ from simulation scenario III. Similar figures can be examined for any particular data set. Unfortunately, overall increases in bias entirely offset decreases in variance so that the mean squared error is actually slightly higher when using $t_{\mathrm{opt}}$ as compared to $t_{\max}$, at least in these limited situations. Note that this persists even with a much larger sample size. In large part, this is due to the tendency of the Cox and Greenwood variance estimators to underestimate the true variability so that $t_{\mathrm{opt}}$ is selected to allow for too much bias. This phenomenon is likely to ameliorate somewhat if the bootstrap variance is used instead.
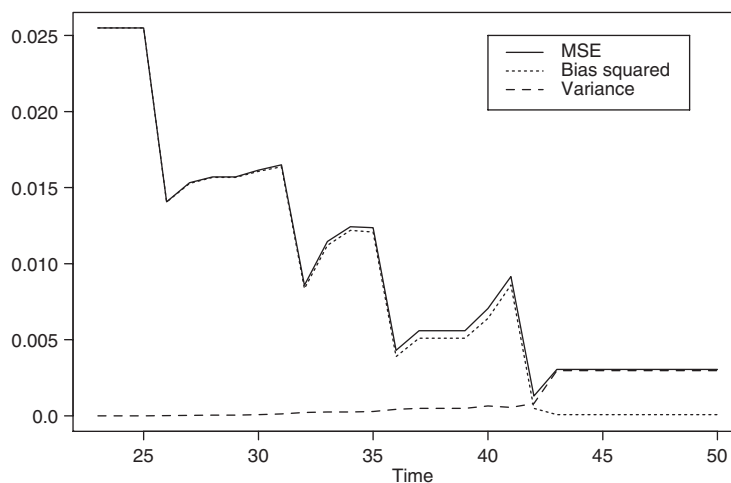
Figure 3. Empirical estimates of the bias squared, variance and mean squared error of $CFR_b$ evaluated at times $t$ lower than $t_{max}$ for simulation scenario III ($n = 100$).

For $n = 100$, the simulation means of the four variance estimators are all reasonably close to their actual variance, with somewhat better performance for $CFR_b$ than for $CFR_a$. For $CFR_b$ evaluated at $t_{opt}$, the Cox approximation does not perform as well with a tendency to substantially underestimate variability in scenarios II and III as compared to the Greenwood estimator.

For $n = 1500$, the situation is now qualitatively similar although here the influence curve method substantially overstates the variability in all three scenarios, presumably because estimates are made further out in the tail of $F_1$ and $F_2$ as evidenced by the lower bias throughout with the higher sample size. The bootstrap method continues to exhibit reasonable performance, outperforming both the Greenwood estimator and the Cox approximation which, while much better than the influence curve calculation, significantly underestimate the true variance of estimators. As with the smaller sample size, the underestimation associated with the Cox approximation is even greater in estimating the variance of $CFR_b$ evaluated at $t_{opt}$. The bias associated with the heavier censoring of scenarios I and II yield much smaller confidence interval coverage than is desirable. Only in scenario I are the coverages broadly acceptable.

Overall, the results support preference for the estimator $CFR_b$ here, with no improvement in accuracy when evaluating it at $t_{opt}$ as against $t_{max}$. For variance estimation, the bootstrap method is most effective; plug-in evaluation of the influence curve variance cannot be recommended in general. The Cox approximation is surprisingly competitive with the Greenwood formula although both exhibit a tendency to underestimate the true variability, particularly when evaluated at $t_{opt}$.

## 4. APPLICATION TO SARS DATA

Severe Acute Respiratory Syndrome (SARS) caused by a previously unknown coronavirus, infected over 8000 people worldwide during 2003. The source of infection of SARS in Hong Kong was

traced to an infected Guangzhou professor who arrived in Hong Kong in late February 2003. The analyses here are based on subsamples of the complete record of 1755 SARS cases in Hong Kong as defined by the WHO clinical case definition. Further epidemiological details for these cases are given in Donnelly *et al.* [2] and Leung *et al.* [18].

Since the data here contain only those admitted to hospital, the date of the latter event is used as the time origin rather than the date of infection. The original data included 124 patients who were in hospital prior to onset of SARS and who therefore had other, possibly serious, conditions that could potentially bias upwards estimates of the CFR; hence it is sensible to exclude these patients when estimating a CFR intended to be applicable to the general population. The date of recovery refers to the final date of discharge from a health-care facility (not just the acute care hospital of first admission)—no recovered individuals subsequently died of SARS-related causes. A further 25 cases were not used because either their final outcome or discharge date was unknown. To illustrate the estimators and their associated variability, we analyse the data as they would have been observed at seven different time points, weekly starting with 2 April 2003, during which period the epidemic almost doubled in size.

Table III provides the estimates $CFR_a$ and $CFR_b$, along with various variability estimates, at each of the seven times noted (as before we used 200 bootstrap samples with replacement). The level of censoring at the various analysis times is provided. As background, the mean time to outcome (death or recovery) was a little more than three weeks. The epidemic of infections peaked in late March, but given the lag time between infection and outcome, 84 per cent of existing case outcomes remained unknown at the beginning of April, reflecting very heavy censoring for the earliest analyses of Table III. Given the simulation results of Section 3, we thus expect CFR estimates in the first four weeks of April to exhibit large bias.

Given the simulation results of Section 3, we expect CFR estimates in the first four weeks of April to exhibit large bias, and this is clearly reflected in the results of Table III. We note that the estimate $CFR_b$ was relatively stable after late April, with a value of 14.2 per cent for the analysis of 14 May. The alternative estimator $CFR_a$ gave slightly lower values. With the simulation evidence of Section 3 for similar shapes for $F_1$ and $F_2$, the estimator $CFR_b$ is preferred here. We note that the case fatality ratio 14.2 per cent is lower than the final reported value of 17.2 per cent for Hong Kong, the discrepancy largely due to the exclusion of the 124 patients who contracted SARS after admission to hospital, a group that displayed a much greater case fatality ratio than the general population.

As anticipated from the simulation results of Section 3, the variance estimates from the Greenwood formula, the Cox approximation, or the bootstrap appear much more reliable than those based on simple estimation of the influence curve of $CFR_b$, although the latter gives acceptable results for the last analysis, at least for $CFR_b$. There is little to choose here between the Cox approximation and the bootstrap—both give very similar results. Finally, bias quickly dominates gains in precision in choosing smaller values of time at which to estimate the CFR so that here $CFR_b(t_{opt})$ is close to $CFR_b$. Finally, the simulations of Section 3 suggest that the Greenwood formula is more reliable as an estimator of the associated $\widehat{var}(CFR_b(t_{opt}))$ than the Cox approximation, although it still likely underestimates the true variability.

In summary, it is clear that CFR estimators have substantial bias when the level of censoring exceeds 40 per cent and are therefore unreliable. Overall, and for the later analyses in particular, the bootstrap variance estimator seems to be the preferred variance estimator for either approach.

Table III. Various estimates of the case fatality ratio for SARS at different time points during the 2003 Hong Kong epidemic.

| | 2 April | 9 April | 16 April | 23 April | 30 April | 7 May | 14 May |
|---|---|---|---|---|---|---|---|
| Sample size | 929 | 1198 | 1352 | 1466 | 1524 | 1558 | 1583 |
| % of Observations censored | 83.6 | 74.5 | 57.6 | 44.7 | 30.2 | 20.3 | 13.9 |
| $CFR_a$ $(t_{max})$ | 0.0743 | 0.0640 | 0.0832 | 0.1068 | 0.1220 | 0.1347 | 0.1370 |
| $\widehat{var}(CFR_a)$ (IC) | 0.08801 | 0.03500 | 0.02728 | 0.01490 | 0.00518 | 0.00206 | 0.00074 |
| $\widehat{var}(CFR_a)$ (Cox) | 0.00064 | 0.00024 | 0.00020 | 0.00022 | 0.00015 | 0.00014 | 0.00013 |
| $\widehat{var}(CFR_a)$ (Greenwood) | 0.00020 | 0.00013 | 0.00013 | 0.00014 | 0.00011 | 0.00010 | 0.00009 |
| $CFR_b$ $(t_{max})$ | 0.0743 | 0.0743 | 0.0971 | 0.1220 | 0.1334 | 0.1422 | 0.1419 |
| $\widehat{var}(CFR_b)$ (IC) | 0.00310 | 0.00234 | 0.00033 | 0.00087 | 0.00043 | 0.00026 | 0.00015 |
| $\widehat{var}(CFR_b)$ (Cox) | 0.00078 | 0.00031 | 0.00024 | 0.00023 | 0.00015 | 0.00012 | 0.00011 |
| $\widehat{var}(CFR_b)$ (Greenwood) | 0.00026 | 0.00027 | 0.00028 | 0.00026 | 0.00020 | 0.00019 | 0.00017 |
| $\widehat{var}(CFR_b)$ (bootstrap) | 0.00079 | 0.00034 | 0.00026 | 0.00021 | 0.00013 | 0.00011 | 0.00010 |
| $CFR_b$ $(t_{opt})$ (Cox) | 0.0664 | 0.0726 | 0.0893 | 0.1219 | 0.1310 | 0.1419 | 0.1401 |
| $\widehat{var}(CFR_b (t_{opt}))$ (Cox) | 0.00040 | 0.00029 | 0.00019 | 0.00020 | 0.00013 | 0.00012 | 0.00010 |
| $t_{opt}$ (days) | 18 | 24 | 18 | 19 | 33 | 57 | 52 |
| $CFR_b (t_{opt})$ (Greenwood) | 0.0668 | 0.0678 | 0.0934 | 0.1163 | 0.1334 | 0.1419 | 0.1400 |
| $\widehat{var}(CFR_b (t_{opt}))$ (Greenwood) | 0.00019 | 0.00009 | 0.00014 | 0.00013 | 0.00015 | 0.00015 | 0.00014 |
| $t_{opt}$ (days) | 18 | 17 | 23 | 26 | 32 | 40 | 42 |

## 5. DISCUSSION

For the Hong Kong SARS data we have argued in favour of the estimator $\mathrm{CFR}_b$ over $\mathrm{CFR}_a$. While this preference remains reasonable when the cumulative sub-distributions of death and recovery increase at a similar rate over time, this need not always be so. For example, consider an alternative scenario where those patients, at risk of dying, tend to die very quickly after infection, whereas it typically takes somewhat longer for recovery. In this case, the death sub-distribution, $F_1$ reaches its asymptote at a much smaller $t$, than $F_2$. Then, for low values of $t$ (and therefore early in an epidemic) $\mathrm{CFR}_b$ suffers from much greater bias than $\mathrm{CFR}_a$ as illustrated in Figure 4.

Estimation of the relevant influence curve of a case fatality ratio estimator as the basis of a variance calculation is unreliable when there is substantial censoring. We note that this does not contradict the expression of the asymptotic variance in terms of the influence curve, only that simple plug-in estimators will not perform well unless $\bar{G}$ is away from zero. Alternative estimators of the same asymptotic variance have much better performance. This is not surprising in light of the standard Kaplan–Meier estimator where Greenwood's formula is much more stable than plug-in estimates of the influence curve under heavy censoring. Both the Greenwood estimator and the Cox approximation are plausible alternatives in the situations considered here, although it remains to be proved that the latter approximation is always asymptotically correct. However, the bootstrap variance estimate is even more attractive, although more computationally intensive. It appears as if the bootstrap technique is picking up second-order effects in variance estimation that are important in estimation in the tails of the sub-distribution functions and that are sometimes missed by first order asymptotic estimators.

Evidence has suggested that the case fatality ratio for SARS varies with other patient cofactors, principally age, where the elderly suffer from far greater case fatality [2, 18]. The procedures studied here can be immediately applied to subgroups of interest. Extending our results to regression analyses that allow the case fatality ratio to vary continuously with an interval-scaled explanatory
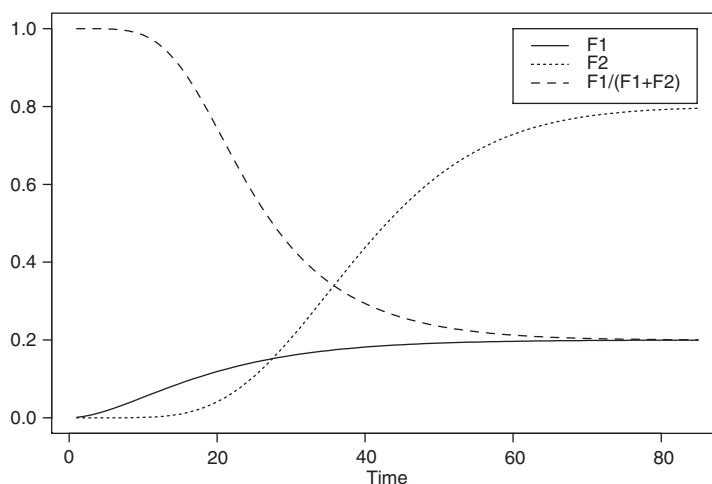


Figure 4. Description of the sub-distribution functions $F_1$, $F_2$ and ratio $F_1/(F_1 + F_2)$ for scenario where times to death are considerably smaller than to recovery.

variable is a topic that is currently being pursued using standard regression models for competing risks data.

## APPENDIX A: INFLUENCE CURVES

With independent right censored observation of competing risks, the full data $X = (T, J, C)$ is subject to an observation process that is coarsened at random (CAR); that is, given the full data $X$, the coarsening mechanism only depends on the observed data (Chapter 1 [19]). From van der Laan and Robins [19] the efficient influence curve for estimation of $F_1(t)$ is then given by

$$\text{IC}_{\text{eff}}(Y; t) = \frac{(I(Y \leqslant t, J = 1) - F_1(t))\Delta}{\bar{G}(T)}$$

$$+ \int_0^{\tilde{T}} \frac{P(T \leqslant t, J = 1 | T \geqslant u) - F_1(t)}{\bar{G}(u)} \, \mathrm{d}M_G(u) \tag{A1}$$

where $\bar{G}(T) = 1 - G(T)$, and

$$\mathrm{d}M_G(u) = I(\tilde{T} \in \mathrm{d}u, \Delta = 0) - I(\tilde{T} \geqslant u)\Lambda_G(\mathrm{d}u)$$

with $\Lambda_G$ denoting the cumulative hazard function associated with $G$.

It remains to establish that $\text{IC}_1 \equiv \text{IC}_{\text{eff}}$. However, this immediately follows from the fact that $\hat{F}_1$ is the non-parametric maximum likelihood estimator under CAR and is also asymptotically linear [20].

For calculation of the influence curve for $\text{CFR}_b$, simple algebra shows that

$$\frac{\hat{F}_1(t)}{\hat{F}_1(t) + \hat{F}_2(t)} - \frac{F_1(t)}{F_1(t) + F_2(t)} = \frac{1}{F(t)}(\hat{F}_1(t) - F_1(t)) - \frac{\hat{F}_1(t)}{\hat{F}(t)F(t)}(\hat{F}(t) - F(t))$$

where $F(t) = F_1(t) + F_2(t)$ and $\hat{F}(t) = \hat{F}_1(t) + \hat{F}_2(t)$.

Thus, for fixed $t$, the influence curve, IC, for $\hat{F}_1(t)/(\hat{F}_1(t) + \hat{F}_2(t))$ as an estimator of $F_1(t)/(F_1(t) + F_2(t))$ depends straightforwardly on the influence curves of $F_1(t)$ and $F_2(t)$

$$\text{IC}(Y_i; t) = \frac{1}{F(t)}\text{IC}_1(Y_i; t) - \frac{F_1(t)}{[F(t)]^2}(\text{IC}_1 + \text{IC}_2)(Y_i; t) \tag{A2}$$

## APPENDIX B: GREENWOOD'S FORMULA IN THE COMPETING RISKS SETTING

First, we introduce some necessary notation. For $j = 1, 2$, let $F_{0j}(s, t)$ denote the probability that an individual alive at time $s$ will have failed due to cause $J = j$ by time $t$. Similarly, let $F_{00}(s, t)$ be the probability that an individual alive at time $s$ is still alive at time $t$. (Note that, in the SARS example, alive refers to being infected, and failure refers to either recovery or death.) In particular, $F_{0j}(0, t) = F_j(t)$ and $F_{00}(0, t) = S(t)$. Then, using non-parametric maximum likelihood, we have

the estimators

$$\hat{F}_{00}(s,t) = \prod_{s<t_i\leqslant t}\left(1-\frac{d_i}{n_i}\right)$$

where $d_i = d_{i1}+d_{i2}$, and

$$\hat{F}_{0j}(s,t) = \sum_{s<t_i\leqslant t}\frac{d_{ij}}{n_i}\hat{F}_{00}(s,t_i^-)$$

Then, following Andersen *et al.* [12], a consistent estimator of the covariance of the estimators $\hat{F}_1(t)$ and $\hat{F}_2(t)$ is given by

$$\widehat{\mathrm{Cov}}(\hat{F}_1(t),\hat{F}_2(t)) = -\sum_{t_i\leqslant t}\hat{S}^2(t_i^-)\{1-\hat{F}_{01}(t_i,t)\}\hat{F}_{02}(t_i,t)\times\frac{(n_i-1)}{n_j^3}d_{i1}$$

$$-\sum_{t_i\leqslant t}\hat{S}^2(t_i^-)\hat{F}_{01}(t_i,t)(1-\hat{F}_{02}(t_i,t))\times\frac{(n_i-1)}{n_i^3}d_{i2} \qquad (B1)$$

Further, a Greenwood estimate of the asymptotic variance of $\hat{F}_j(t)$ is given by

$$\widehat{\mathrm{Var}}(\hat{F}_j(t)) = \sum_{t_i\leqslant t}[\hat{S}(t_i^-)\hat{F}_{0j}(t_i,t)]^2\times\frac{(n_i-1)}{n_i^3}(d_{i1}+d_{i2})$$

$$+\sum_{t_i\leqslant t}\hat{S}^2(t_i^-)[1-2\hat{F}_{0j}(t_i,t)]\times\frac{(n_i-1)}{n_i^3}d_{ij} \qquad (B2)$$

## APPENDIX C: THE COX APPROXIMATION FOR $\sigma^2$

Noting that $\log\hat{S}(t) = \sum_{t_i<t}\log(1-\hat{h}_i)$, where $\hat{h}_i\equiv d_i/n_i$, it follows from the delta method that $\widehat{\mathrm{cov}}(\hat{S}(s),\hat{S}(t)) = \hat{S}(s)\hat{S}(t)\sum_{t_i<s}d_i/n_i(n_i-d_i)$ for $s\leqslant t$ estimates the asymptotic covariance of the Kaplan–Meier estimator at two times $s$ and $t$. Cox suggested ignoring the covariance between $\hat{h}_{ij}=d_{ij}/n_i$ and $\hat{S}(t_k^-)$ for any $i,j,k$, so that we can again use the delta method on (1) to yield

$$\widehat{\mathrm{var}}(\hat{F}_1(t)) = \mathbf{h}_1\Omega\mathbf{h}_1^{\mathrm{T}}+\sum_{t_i\leqslant t}\hat{S}(t_i^-)^2\frac{d_i(n_i-d_i)}{n_i^3} \qquad (C1)$$

where $\mathbf{h}_1 = (\hat{h}_{11},\ldots,\hat{h}_{k1})$ with $t_k$ being the largest event time $\leqslant t$, and $\Omega$ the variance–covariance matrix of $\mathbf{S}_k\equiv(\hat{S}(t_1^-),\ldots,\hat{S}(t_k^-))$. An analogous approximation obtains for $\widehat{\mathrm{var}}(\hat{F}_2(t))$. We note that this provides a simple approximation to the variance of CFR$_a$, although, in this case, the estimator (C1) does not reduce to a standard multinomial estimator although the ignored co-variance terms are of smaller order. (In fact, putting in the correct covariance terms for all the relevant $\hat{h}_{ij}$ and $\hat{S}(t_k^-)$ provides an alternative method for computation of a Greenwood variance formula.)

By the same approach, we also obtain the approximation

$$\widehat{\text{cov}}(\hat{F}_1(t), \hat{F}_2(t)) = \mathbf{h}_1 \mathbf{\Omega} \mathbf{h}_2^{\text{T}} + \mathbf{S}_k \mathbf{\Psi} \mathbf{S}_k^{\text{T}} \tag{C2}$$

where $\Psi$ is an estimate of the asymptotic variance–covariance matrix of the vectors $\mathbf{h}_1$ and $\mathbf{h}_2$, thus given by a diagonal matrix with $(i, i)$ entry equal to $-d_{i1}d_{i2}/n_i{}^2$. A final application of the delta method, using (2), gives the final formula for the estimator of the asymptotic variance of $\text{var}(\hat{F}_1(t)/(\hat{F}_1(t) + \hat{F}_2(t)))$. We note that Ghani *et al.* [4] ignore certain 'second-order' terms so that they take the summand in the second term of (C1) to be $\hat{S}(t_i{}^-)^2(d_i/n_i{}^2)$ and $\Psi = 0$; we use this version throughout the article.

## ACKNOWLEDGEMENTS

## REFERENCES

1. Gaynor JJ, Feuer EJ, Tan CC, Wu DH, Little CR, Straus DJ, Clarkson BD, Brennan MF. On the use of cause-specific failure and conditional failure probabilities: examples from clinical oncology data. *Journal of the American Statistical Association* 1993; **88**:400–409.
2. Donnelly CA, Ghani AC, Leung GM, Hedley AJ, Fraser C, Riley S, Abu-Raddad LJ, Ho L-M, Thach T-Q, Chau P, Chan K-P, Lam T-H, Tse L-Y, Tsang T, Liu S-H, Kong JHB, Lau EMC, Ferguson NM, Anderson RM. Epidemiological determinants of spread of causal agent of severe acute respiratory syndrome in Hong Kong. *The Lancet* 2003; **361**:1761–1766.
3. Galvani AP, Lei X, Jewell NP. Severe Acute Respiratory Syndrome: temporal stability and geographic variation in cumulative case fatality rates and average doubling times. *Emerging Infectious Diseases* 2003; **9**:991–994.
4. Ghani AC, Donnelly CA, Cox DR, Griffin JT, Fraser C, Lam TH, Ho LM, Chan WS, Anderson RM, Hedley AJ, Leung GM. Methods for estimating the case fatality rate for a novel, emerging infectious disease. *American Journal of Epidemiology* 2005; **162**:479–486.
5. Kalbfleisch JD, Prentice RL. *The Statistical Analysis of Failure Time Data* (2nd edn). Wiley: New York, 2002.
6. Koul H, Susarla V, Van Ryzin J. Regression analysis with randomly right censored data. *Annals of Statistics* 1981; **9**:1276–1288.
7. Robins JM, Rotnitzky A. Recovery of information and adjustment for dependent censoring using surrogate markers. In *AIDS Epidemiology*: *Methodological Issues*, Jewell NP, Dietz K, Farewell V (eds). Birkhauser: Boston, 1992; 297–331.
8. Aalen OO, Johansen S. An empirical transition matrix for nonhomogeneous Markov chains based on censored observations. *Scandinavian Journal of Statistics* 1978; **5**:141–150.
9. Aalen OO. Nonparametric estimation of partial transition probabilities in multiple decrement models. *Annals of Statistics* 1978; **6**:534–545.
10. Gooley TA, Leisenring W, Crowley J, Storer BE. Estimation of failure probabilities in the presence of competing risks: new representations of old estimators. *Statistics in Medicine* 1999; **18**:695–706.
11. Reid N. Influence functions. In *Encyclopedia of Statistical Sciences*, vol. 4, Kotz S, Johnson NL (eds). Wiley: New York, 1983; 117–120.
12. Andersen PK, Borgan O, Gill RD, Keiding N. *Statistical Models Based on Counting Processes*. Springer: New York, 1993.
13. Dinse GE, Larson MG. A note on semi-Markov models for partially censored data. *Biometrika* 1986; **73**:379–386.
14. Pepe MS. Inference for events with dependent risks in multiple endpoint studies. *Journal of the American Statistical Association* 1991; **86**:770–778.
15. Lin DY. Non-parametric inference for cumulative incidence functions in competing risks studies. *Statistics in Medicine* 1997; **16**:901–910.
16. Choudhury JB. Non-parametric confidence interval estimation for competing risks analysis: application to contraceptive data. *Statistics in Medicine* 2002; **21**:1129–1144.

17. Shao J, Tu D. *The Jackknife and Bootstrap*. Springer: New York, 1995.
18. Leung GM, Hedley AJ, Ho L-M, Chau P, Wong IOL, Thach TQ, Ghani AC, Donnelly CA, Fraser C, Riley S, Ferguson NM, Anderson RM, Tsang T, Leung P-Y, Wong V, Chann JCK, Tsui E, Lo S-V, Lam T-H. The epidemiology of severe acute respiratory syndrome in the 2003 Hong Kong epidemic: an analysis of all 1755 patients. *Annals of Internal Medicine* 2004; **141**:662–673.
19. van der Laan MJ, Robins JM. *Unified Methods for Censored Longitudinal Data and Causality*. Springer: New York, 2003.
20. Gill RD, van der Laan MJ, Robins JM. Coarsening at random: characterizations, conjectures, counter-examples. In *Proceedings of First Seattle Conference in Biostatistics*, Lin D-Y (ed.). Springer: Berlin, 1997; 255–294.