

Covariance formulas via marginal martingales

J. A. Wellner¹

Department of Statistics, GN-22
 University of Washington
 Seattle, WA 98195, USA

Prentice and Cai recently introduced and studied the function C defined as the covariance function of the two marginal counting process martingales of a pair of dependent survival times (T_1, T_2) . They show that the function C together with the marginal distributions determines the joint survival function F of (T_1, T_2) . In this note we show how the key characterizing equation of Prentice and Cai yields a formula for the covariance of T_1 and T_2 in terms of the marginal mean residual life functions and C . The resulting formula generalizes a formula for the variance of a one-dimensional random variable T due to Pyke (1965). We also explore several generalizations of the covariance formula, and obtain a valid k -dimensional version of the Prentice and Cai formula.

Key Words & Phrases: bivariate distribution, covariance, Hoeffding's formula, marginal martingales, mean residual life functions.

Let T be a non-negative random variable with distribution function $F: F(t) = P(T \leq t)$. Then the survival function is $\bar{F}(t) = 1 - F(t) = P(T > t)$, and the cumulative hazard function is $A(t) = \int_{[0,t]} (1 - F_-)^{-1} dF$. Let $N(t) = 1_{[T \leq t]}$ be the corresponding elementary counting process, and let M be the counting process martingale defined by

$$M(t) = N(t) - \int_0^t 1_{[T \geq s]} dA(s), \quad t \geq 0. \quad (1)$$

Define V for $t \geq 0$ by

$$\begin{aligned} V(t) &\equiv \text{Var}(M(t)) = E \langle M \rangle (t) = \\ &= E \int_0^t 1_{[T \geq s]} (1 - \Delta A(s)) dA(s) = \int_0^t (1 - \Delta A(s)) dF(s) \end{aligned} \quad (2)$$

where $\Delta A(t) = A(t) - A(t-)$. Note that $V = F$ when F , and hence also A , is continuous. Assuming that $ET < \infty$, the mean residual life function e of T is defined by

$$e(t) = E(T - t | T > t) = \frac{\int_t^\infty \bar{F}(s) ds}{\bar{F}(t)}, \quad t \geq 0. \quad (3)$$

¹ Research supported in part by: National Science Foundation Grant DMS-9108409 and NWO Grant B 61-238.

When the distribution function F of T is continuous and $ET^2 < \infty$, PYKE (1965) discovered the following formula giving the variance of T in terms of the mean residual life function e and $V(t) \equiv \text{Var}(M(t)) = F(t)$:

$$\text{Var}(T) = \int e^2(t) dV(t). \quad (4)$$

Note that $V \neq F$ when F is discontinuous. The formula (4) was obtained for an arbitrary distribution function F (in a somewhat different form) by HALL and WELLNER (1981). SHORACK and WELLNER (1986), page 283, generalized it to an arbitrary distribution function F and an arbitrary square integrable function a of T :

$$\text{Var}[a(T)] = E[e_a^2(T)(1 - \Delta A(T))] = \int_0^\infty e_a^2(t) dV(t) \quad (5)$$

where

$$e_a(t) = E(a(T) - a(t) | T > t) = -Ra(t). \quad (6)$$

Also see e.g. RITOV and WELLNER (1988), remark 2.2, page 201, or EFRON and JOHNSTONE (1990), for connections with the R and L operators arising in survival analysis ($e_a = -Ra$ and $L = R^T = R^{-1}$ is a martingale operator).

Now let (T_1, T_2) be a pair of survival times with joint survival function \bar{F} : $\bar{F}(t_1, t_2) = P(T_1 > t_1, T_2 > t_2)$. Let $N_1(t_1) = 1[T_1 \leq t_1]$ and $N_2(t_2) = 1[T_2 \leq t_2]$ be the two marginal counting processes, and let M_1, M_2 be the marginal counting process martingales defined by

$$M_i(t_i) = N_i(t_i) - \int_0^{t_i} 1_{[T_i \geq s]} dA_i(s), \quad t_i \geq 0, \quad i = 1, 2. \quad (7)$$

PRENTICE and CAI have recently introduced and studied the function C defined by

$$C(t_1, t_2) \equiv E\{M_1(t_1)M_2(t_2)\}, \quad t_i \geq 0, \quad i = 1, 2. \quad (8)$$

They show that C , together with the marginal survival functions $\bar{F}_i(t_i) = P(T_i > t_i)$, $i = 1, 2$, determine the joint survival function \bar{F} of T_1, T_2 by the following formula:

$$\bar{F}(t_1, t_2) = \bar{F}_1(t_1)\bar{F}_2(t_2) \left[1 + \int_0^{t_1} \int_0^{t_2} \frac{1}{\bar{F}_1(s_1)\bar{F}_2(s_2)} C(ds_1, ds_2) \right]; \quad (9)$$

see equation (8) of PRENTICE and CAI (1992a) or (1992b). This interesting representation of \bar{F} in terms of C and the marginal distributions F_1 and F_2 has applications to a wide range of bivariate estimation problems, including the bivariate random right censoring model. By use of (8) one can obtain estimators of \bar{F} ; see e.g. PRENTICE and CAI (1992a), (1992b), and GILL, VAN DER LAAN and WELLNER (1994).

PRENTICE and CAI (1992a) also derive a k -dimensional extension of (9); see their formula (23), page 508. However, as pointed out by GILL (1992), their k -dimensional formula is incorrect: see PRENTICE and CAI (1993).

The object of this note is to show how (9) yields a formula analogous to (4) for

$\text{Cov}(T_1, T_2)$ in terms of C and the marginal mean residual life functions e_i , $i = 1, 2$, defined for $i = 1, 2$ by

$$e_i(t_i) = E(T_i - t_i | T_i > t_i) = \frac{\int_{t_i}^{\infty} \bar{F}_i(s_i) ds_i}{\bar{F}_i(t_i)}, \quad t_i \geq 0. \quad (10)$$

We will also give a bivariate analogue of (5), and a multivariate generalization thereof. The latter yields, in turn, a correct multivariate generalization of (9).

PROPOSITION 1. *Suppose that (T_1, T_2) have joint survival function \bar{F} , marginal survival and mean residual life functions \bar{F}_i , e_i , $i = 1, 2$, and that C is defined by (8). Then, if $\text{Var}(T_i) < \infty$, $i = 1, 2$,*

$$\text{Cov}(T_1, T_2) = \int_0^{\infty} \int_0^{\infty} e_1(t_1) e_2(t_2) C(dt_1, dt_2). \quad (11)$$

Furthermore, the bivariate analogue of (5) holds: for any measurable functions $a_i(T_i)$, $i = 1, 2$ with $\text{Var}(a_i(T_i)) < \infty$, $i = 1, 2$,

$$\text{Cov}(a_1(T_1), a_2(T_2)) = \int_0^{\infty} \int_0^{\infty} e_{a_1}(t_1) e_{a_2}(t_2) C(dt_1, dt_2), \quad (12)$$

where e_{a_i} , $i = 1, 2$ are given by (6) (with a, T, t replaced by a_i, T_i, t , $i = 1, 2$).

PROOF. Our first proof of (11) and (12) will proceed directly from (9); the second proof will yield (12) directly via the martingale representations of RITOV and WELLNER (1988). By a well-known formula due to Hoeffding (see e.g. LEHMANN, 1966, or SHORACK and WELLNER, 1986, formula (A.8.2), page 862),

$$\text{Cov}(T_1, T_2) = \int_0^{\infty} \int_0^{\infty} (\bar{F}(t_1, t_2) - \bar{F}_1(t_1)\bar{F}_2(t_2)) dt_1 dt_2. \quad (a)$$

Hoeffding's original version of this formula applies to distribution functions, but the current version in terms of survival functions is an elementary consequence.

It then follows immediately, by using PRENTICE and CAI's formula (9) to evaluate the integrand on the right side in (a), and then using Fubini's theorem, that

$$\text{Cov}(T_1, T_2) = \iint \bar{F}_1(t_1) \bar{F}_2(t_2) \int_0^{t_1} \int_0^{t_2} \frac{1}{\bar{F}_1(s_1) \bar{F}_2(s_2)} C(ds_1, ds_2) dt_1 dt_2 \quad (b)$$

$$\begin{aligned} &= \iiint \frac{1_{[s_1 \leq t_1]} 1_{[s_2 \leq t_2]} \bar{F}_1(t_1) \bar{F}_2(t_2)}{\bar{F}_1(s_1) \bar{F}_2(s_2)} C(ds_1, ds_2) dt_1 dt_2 \\ &= \iint \frac{\int_{[s_1 < t_1]} \bar{F}_1(t_1) dt_1}{\bar{F}_1(s_1)} \frac{\int_{[s_2 < t_2]} \bar{F}_2(t_2) dt_2}{\bar{F}_2(s_2)} C(ds_1, ds_2) \\ &= \iint e_1(s_1) e_2(s_2) C(ds_1, ds_2). \quad (c) \end{aligned}$$

The second formula (12) follows by first observing that (9) implies that it holds for the indicator functions $a_i(T_i) = 1_{[T_i > t_i]}$, $i = 1, 2$ for any (t_1, t_2) , since $e_{a_i}(s) = -R_F a_i(s) = -1_{[0, t_i]}(s) \bar{F}(t) / \bar{F}(s)$ if $a_i(s) = 1_{[s > t]}$. By bilinearity of both sides of the equation, this yields (12) by the standard argument of obtaining measurable functions as limits of simple functions.

Here is a second proof of (12) via the martingale representations of RITOV and WELLNER (1988). By (2.8) and (2.7) of RITOV and WELLNER (1988),

$$a_i(T_i) - E a_i(T_i) = L_i \circ R_i a_i(T_i) = \int_0^\infty R_i a_i dM_i, \quad i = 1, 2 \tag{d}$$

where L_i and R_i are the L and R operators corresponding to the marginal distribution functions F_i , $i = 1, 2$ respectively; here M_i , $i = 1, 2$ are the marginal counting process martingales given by (7). Also, from formula (2.4) in RITOV and WELLNER (1988), $R_i a_i = -e_{a_i}$.

The representations given in (d) yield

$$\begin{aligned} \text{Cov}(a_1(T_1), a_2(T_2)) &= E \left(\int_0^\infty R_1 a_1 dM_1 \int_0^\infty R_2 a_2 dM_2 \right) \tag{e} \\ &= \int_0^\infty \int_0^\infty R_1 a_1(t_1) R_2 a_2(t_2) C(dt_1, dt_2) \\ &= \int_0^\infty \int_0^\infty e_{a_1}(t_1) e_{a_2}(t_2) C(dt_1, dt_2) \end{aligned}$$

by application of Fubini's theorem. □

The second proof of (12) given above immediately yields the following multivariate generalization of proposition 1:

PROPOSITION 2. Suppose that $T = (T_1, \dots, T_k)$ has distribution function F on R^{+k} and that $a_i(T_i)$ are measurable functions with $E(a_i^k(T_i)) < \infty$, $i = 1, \dots, k$. Then

$$\begin{aligned} &E \left\{ \prod_{i=1}^k (a_i(T_i) - E a_i(T_i)) \right\} \\ &= (-1)^k \int_0^\infty \dots \int_0^\infty \left\{ \prod_{i=1}^k e_{a_i}(t_i) \right\} C(dt_1, \dots, dt_k) \tag{13} \end{aligned}$$

where

$$C(t_1, \dots, t_k) \equiv E \prod_{i=1}^k M_i(t_i). \tag{14}$$

PROOF. Use the representation in (d), now for $i = 1, \dots, k$, and proceed exactly as in (e). □

When specialized to the indicator functions $a_i(T_i) = 1_{[T_i > t_i]}$, $i = 1, \dots, k$, (13) yields one possible k -variate generalization of (9) as an immediate corollary:

COROLLARY.

$$\begin{aligned} \bar{F}(t_1, \dots, t_k) &= \psi_k(t_1, \dots, t_k) \\ &+ \bar{F}_1(t_1) \dots \bar{F}_k(t_k) \int_0^{t_1} \dots \int_0^{t_k} \frac{1}{\bar{F}_1(s_1) \dots \bar{F}_k(s_k)} C(ds_1, \dots, ds_k) \end{aligned} \quad (15)$$

where

$$\psi_k(t_1, \dots, t_k) = E \left\{ \sum_{j=1}^k \left(\prod_{i=1}^{j-1} 1_{[T_i > t_i]} \right) \bar{F}_j(t_j) \prod_{i=j+1}^k (1_{[T_i > t_i]} - \bar{F}_i(t_i)) \right\} \quad (16)$$

is a function of the $k-1$ and lower dimensional marginal survival functions of F .

PROOF. First note that for $a_i(T_i) = 1_{[T_i > t_i]}$,

$$e_{a_i}(s_i) = -R_{F_i} a_i(s) = -1_{[0, t_i]}(s_i) \frac{\bar{F}_i(t_i)}{\bar{F}_i(s_i)}, \quad i = 1, \dots, k.$$

Hence (13) yields, with this choice of the functions a_i ,

$$\begin{aligned} &E \prod_{i=1}^k (1_{[T_i > t_i]} - \bar{F}_i(t_i)) \\ &= \bar{F}_1(t_1) \dots \bar{F}_k(t_k) \int_0^{t_1} \dots \int_0^{t_k} \frac{1}{\bar{F}_1(s_1) \dots \bar{F}_k(s_k)} C(ds_1, ds_k). \end{aligned} \quad (a)$$

Rewrite the left side of (a) as

$$\bar{F}(t_1, \dots, t_k) + E \prod_{i=1}^k (1_{[T_i > t_i]} - \bar{F}_i(t_i)) - \bar{F}(t_1, \dots, t_k).$$

Now use the formula

$$\prod_{j=1}^k a_j - \prod_{j=1}^k b_j = \sum_{j=1}^k \left(\prod_{i=1}^{j-1} a_i \right) (a_j - b_j) \left(\prod_{i=j+1}^k b_i \right)$$

to compute

$$\begin{aligned} &E \prod_{i=1}^k (1_{[T_i > t_i]} - \bar{F}_i(t_i)) - \bar{F}(t_1, \dots, t_k) \\ &= -E \left\{ \prod_{i=1}^k 1_{[T_i > t_i]} - \prod_{i=1}^k (1_{[T_i > t_i]} - \bar{F}_i(t_i)) \right\} \\ &= -E \left\{ \sum_{j=1}^k \left(\prod_{i=1}^{j-1} 1_{[T_i > t_i]} \right) \bar{F}_j(t_j) \prod_{i=j+1}^k (1_{[T_i > t_i]} - \bar{F}_i(t_i)) \right\}. \end{aligned}$$

This yields (15). □

Here are three simple examples of (11).

EXAMPLE 1. Suppose that F is Gumbel's bivariate exponential distribution, as considered by PRENTICE and CAI (1992). Thus

$$\bar{F}(t_1, t_2) = e^{-(t_1+t_2)}(1 + \theta(1 - e^{-t_1})(1 - e^{-t_2}))$$

and, as noted by PRENTICE and CAI,

$$C(t_1, t_2) = \frac{\theta}{4} (1 - e^{-2t_1})(1 - e^{-2t_2})$$

$$C(dt_1, dt_2) = \theta e^{-2(t_1+t_2)} dt_1 dt_2.$$

Since the marginal distributions are both exponential, $e_1(t_1) = e_2(t_2) = 1$, and we find that

$$\int_0^\infty \int_0^\infty e_1(t_1)e_2(t_2)C(dt_1, dt_2) = \int_0^\infty \int_0^\infty 1C(dt_1, dt_2) = \frac{\theta}{4}.$$

On the other hand, by direct calculation the density is

$$f(t_1, t_2) = e^{-(t_1+t_2)}(1 + \theta(1 - 2e^{-t_1})(1 - 2e^{-t_2}))1_{[0, \infty)}(t_1)1_{[0, \infty)}(t_2),$$

and this yields

$$E(T_1 T_2) = 1 + \frac{\theta}{4}$$

and hence $\text{Cov}(T_1, T_2) = \theta/4$.

EXAMPLE 2. Suppose that F is the Morgenstern copula function (distribution function on $[0, 1]^2$ with uniform marginals); thus the distribution function is

$$F(u, v) = uv(1 + \theta(1 - u)(1 - v))$$

where $|\theta| < 1$; the density function is

$$f(u, v) = 1 + \theta(1 - 2u)(1 - 2v).$$

It is easy to compute $\text{Cov}(U, V) = \frac{\theta}{36}$. The marginal distributions are uniform $(0, 1)$, and

$$e_1(u) = e_2(u) = \frac{1}{2}(1 - u).$$

Straightforward calculation yields $C(du, dv) = \theta(1 - u)(1 - v) du dv$, and hence

$$\begin{aligned} & \int_0^1 \int_0^1 e_1(u)e_2(v)C(du, dv) \\ &= \frac{1}{4} \int_0^1 \int_0^1 (1 - u)(1 - v)\theta(1 - u)(1 - v) du dv \\ &= \frac{\theta}{4} \left(\int_0^1 (1 - u)^2 du \right)^2 = \frac{\theta}{36}. \end{aligned}$$

EXAMPLE 3. Suppose that F puts mass $1/4$ at each of the four points $(1, 2)$, $(2, 1)$, $(2, 2)$, and $(3, 3)$. Then $ET_1 = 2 = ET_2$, $\text{Var}(T_1) = 1/2 = \text{Var}(T_2)$, $E(T_1 T_2) = 17/4$, and $\text{Cov}(T_1, T_2) = 1/4$.

In this case the two mean residual life functions are given by

$$e_i(t) = (2-t)1_{[0,1)}(t) + (7/3-t)1_{[1,2)}(t) + (3-t)1_{[2,3)}(t), \quad i = 1, 2.$$

so that $e_i(1) = 1$, $e_i(2) = 4/3$, $i = 1, 2$. The function C has jumps at the nine points on the lattice determined by the marginal distributions, and we calculate $C(\Delta 1, \Delta 1) = -1/16$, $C(\Delta 1, \Delta 2) = 1/12 = C(\Delta 1, \Delta 2)$, $C(\Delta 2, \Delta 1) = 5/36$.

Note that the other 5 points (with at least one coordinate equal to 3) do not affect the outcome of our calculation since $e_1(3) = e_2(3) = 0$. Thus we calculate

$$\begin{aligned} & \int_0^\infty \int_0^\infty e_1(t_1) e_2(t_2) C(dt_1, dt_2) \\ &= \frac{4}{3} \frac{4}{3} \left(-\frac{1}{16} \right) + \frac{4}{3} \cdot 1 \cdot \frac{1}{12} + \frac{4}{3} \cdot 1 \cdot \frac{1}{12} + 1 \cdot 1 \cdot \frac{5}{36} \\ &= \frac{1}{4} = \text{Cov}(T_1, T_2). \end{aligned}$$

Acknowledgements

I owe thanks to BRUCE TURNBULL for suggesting the relevance of Hoeffding's formula, RICHARD GILL for help with example 3, and a referee for help in finding the correct version of formula (16).

References

- EFRON, B. and I. JOHNSTONE (1990), Fisher's information in terms of the hazard rate, *Annals of Statistics* 18, 38-62.
- GILL, R. D., *Personal communication*, 1992.
- GILL, R. D., M. J. VAN DER LAAN and J. A. WELLNER (1994), Inefficient estimators of the bivariate survival function for three models to appear in: *Annales de l'Institut Henri Poincaré*.
- HALL, W. J. and J. A. WELLNER (1981), Mean residual life, in: M. CSORGO, D. A. LAWSON, J. N. K. RAO and A. K. MD. E. SALEH (eds.) *Statistics and related topics*, 169-184, North Holland, Amsterdam.
- LEHMANN, E. L. (1966), Some concepts of dependence, *Annals of Mathematical Statistics* 37, 1137-1153.
- PRENTICE, R. L. and J. CAI (1992a), Covariance and survivor function estimation using censored multivariate failure time data, *Biometrika* 79, 495-512.
- PRENTICE, R. L. and J. CAI (1992b), Marginal and conditional models for the analysis of multivariate failure time data, in: J. P. KLEIN and P. K. GOEL (eds.) *Survival Analysis: State of the Art*, 393-406, Kluwer, Dordrecht.
- PRENTICE, R. L. and J. CAI (1993), Amendments and Corrections to: Covariance and survivor function estimation using censored multivariate failure time data, *Biometrika* 80, 711-712.
- PYKE, R. (1965), Spacings, *Journal of the Royal Statistical Society, Series B*, 395-449.
- RITOV, Y. and J. A. WELLNER (1988), Censoring, martingales, and the Cox model, *Contemporary Mathematics* 80, 191-219, American Mathematical Society, Providence.
- SHORACK, G. R. and J. A. WELLNER (1986), *Empirical processes with applications to statistics*, Wiley, New York.

Received: March 1992. Revised: September 1993.