

COMMET 01087

Section II. Systems and programs

## MULCOX: a computer program for the Cox regression analysis of multiple failure time variables

D.Y. Lin

*Department of Biostatistics, Harvard School of Public Health, Boston, MA 02115, U.S.A.*

MULCOX is a user-friendly FORTRAN program for the analysis of regression effects when individual study subjects may experience multiple events or failures. Each marginal distribution of the multivariate failure time variables is formulated by a Cox proportional hazards model. The maximum partial likelihood estimators of the regression parameters in these marginal models are approximately jointly normal. The MULCOX program estimates the marginal models as well as the joint covariance matrix. In addition, it implements several multivariate inference procedures. The program runs on both mainframe computers and microcomputers. The running time is quite acceptable even for large samples. A simple example is provided to illustrate the features of the program.

FORTRAN; Incomplete observations; Multivariate failure times; Proportional hazards; Repeated events; Simultaneous inference; Survival data

### 1. Introduction

Many biomedical studies record the times to two or more distinct events or failures on each subject. The failures may be events of different natures or may be repetitions of the same type of events. The examples of such multivariate failure times include the development of physical symptoms in several major body systems, and the time sequence of asthmatic attacks, infection episodes, tumor diagnoses, or tumor recurrences in individual patients. In these studies, investigators are often interested in assessing the effects of prognostic factors or covariates (e.g., treatment, age and sex) on the multivariate failure time variables.

Since the Cox proportional hazards model [2] is the most commonly used regression technique for analyzing *univariate* failure time data, it is natural to regress each component of the multivariate failure time variables on covariates by a Cox model. In multivariate survival studies, however, it is often desirable to make the statistical inference involving parameters of several failure time variables. For example, we may want to evaluate how the effects of a given covariate vary among failure time variables. Clearly, such a multivariate inference must take into consideration the correlation structure of the parameter estimators in the marginal failure time models. Recently, Wei et al. [4] proved that the parameter estimators of the marginal Cox models are asymptotically jointly normal with a covariance matrix that can be consistently estimated. These authors also proposed various simultaneous inference procedures. In the next Section, we will present the computational

*Correspondence:* D.Y. Lin, Department of Biostatistics, SC-32, University of Washington, Seattle, WA 98195, U.S.A.

methods for this new methodology. The computer program MULCOX was designed to implement these statistical procedures. This program will be described in Section 3 and illustrated with a simple example in Section 4.

## 2. Computational methods

For the  $k$ th failure time variable,  $k = 1, \dots, K$ , let  $V_{ki}$  be the failure time of the  $i$ th subject,  $i = 1, \dots, n$ . In practice, however, we can only observe the bivariate vector  $(X_{ki}, \Delta_{ki})$ , where  $X_{ki} = \min(V_{ki}, C_{ki})$ ,  $C_{ki}$  is the censoring time of the  $i$ th subject with respect to the  $k$ th failure time variable, and  $\Delta_{ki} = 1$  if  $X_{ki} = V_{ki}$  and  $\Delta_{ki} = 0$  otherwise. If  $V_{ki}$  is missing, we let  $C_{ki}$  be 0. This implies that  $X_{ki} = 0$  and  $\Delta_{ki} = 0$  since  $V_{ki}$  is positive. In addition, let  $Z_i = (Z_{1i}, \dots, Z_{pi})'$  denote a  $p \times 1$  vector of covariates for the  $i$ th subject.

The hazard function of the  $k$ th failure time variable for an individual with covariate  $Z$  is assumed to take the form

$$\lambda_k(t; Z) = \lambda_{k0}(t) \exp(\beta'_k Z),$$

where  $\lambda_{k0}(t)$  is an unspecified baseline hazard function, and  $\beta_k = (\beta_{1k}, \dots, \beta_{pk})'$  is the failure-specific regression parameter. The corresponding partial likelihood is

$$L_k(\beta) = \prod_{i=1}^n \left\{ \frac{\exp(\beta'_i Z_i)}{\sum_{j \in \mathcal{R}_k(X_{ki})} \exp(\beta'_j Z_j)} \right\}^{\Delta_{ki}},$$

where  $\mathcal{R}_k(t)$  is the set of labels attached to the subjects at risk just prior to time  $t$  with respect to the  $k$ th failure time variable. Then the maximum partial likelihood estimator  $\hat{\beta}_k$  for  $\beta_k$  is the value of  $\beta$  that maximizes  $L_k(\beta)$ , which is computed by the Gauss-Newton algorithm. The Breslow approximation [1] is used in the case of tied failure times.

Wei et al. [4] showed that, for large  $n$ ,  $\hat{\beta}_T = (\hat{\beta}'_1, \dots, \hat{\beta}'_K)'$  is approximately normal with mean  $\beta_T = (\beta'_1, \dots, \beta'_K)'$  and with joint covariance matrix  $Q$ , say. These authors also provided a consistent estimator  $\hat{Q}$  for  $Q$ .

Before expressing the covariance matrix estimator  $\hat{Q}$ , we need to introduce some notations. First, let  $Y_{ki}(t) = 1$  if  $X_{ki} \geq t$  and  $Y_{ki}(t) = 0$  otherwise. Second, let

$$S_k^{(0)}(\beta, t) = \sum_{i=1}^n Y_{ki}(t) \exp(\beta'_i Z_i),$$

$$S_k^{(1)}(\beta, t) = \sum_{i=1}^n Y_{ki}(t) \exp(\beta'_i Z_i) Z_i,$$

and

$$S_k^{(2)}(\beta, t) = \sum_{i=1}^n Y_{ki}(t) \exp(\beta'_i Z_i) Z_i Z_i'.$$

In addition, let

$$\hat{A}_k(\hat{\beta}_k) = \sum_{i=1}^n \Delta_{ki} \left\{ \frac{S_k^{(2)}(\hat{\beta}_k, X_{ki})}{S_k^{(0)}(\hat{\beta}_k, X_{ki})} - \frac{S_k^{(1)}(\hat{\beta}_k, X_{ki}) S_k^{(1)}(\hat{\beta}_k, X_{ki})'}{S_k^{(0)}(\hat{\beta}_k, X_{ki})^2} \right\},$$

and

$$W_{ki}(\hat{\beta}_k) = \Delta_{ki} \left\{ Z_i - \frac{S_k^{(1)}(\hat{\beta}_k, X_{ki})}{S_k^{(0)}(\hat{\beta}_k, X_{ki})} \right\} - \sum_{j=1}^n \frac{\Delta_{kj} Y_{kj}(X_{ki}) \exp(\hat{\beta}'_k Z_j)}{S_k^{(0)}(\hat{\beta}_k, X_{kj})} \times \left\{ Z_j - \frac{S_k^{(1)}(\hat{\beta}_k, X_{kj})}{S_k^{(0)}(\hat{\beta}_k, X_{kj})} \right\}.$$

Finally, let  $\hat{B}_{kl}(\hat{\beta}_k, \hat{\beta}_l) = \sum_{i=1}^n W_{ki}(\hat{\beta}_k) W_{li}(\hat{\beta}_l)'$ , and  $\hat{D}_{kl}(\hat{\beta}_k, \hat{\beta}_l) = \hat{A}_k^{-1}(\hat{\beta}_k) \hat{B}_{kl}(\hat{\beta}_k, \hat{\beta}_l) \hat{A}_l^{-1}(\hat{\beta}_l)$ . Then the covariance matrix estimator  $\hat{Q}$  is

$$\begin{bmatrix} \hat{D}_{11}(\hat{\beta}_1, \hat{\beta}_1) & \cdots & \hat{D}_{1K}(\hat{\beta}_1, \hat{\beta}_K) \\ \vdots & & \vdots \\ \hat{D}_{K1}(\hat{\beta}_K, \hat{\beta}_1) & \cdots & \hat{D}_{KK}(\hat{\beta}_K, \hat{\beta}_K) \end{bmatrix}.$$

This estimator turns out to be a robust estimator of the covariance matrix of  $\hat{\beta}_T$  [3].

The aforementioned results provide the basis for the simultaneous inference about the effects of covariates on the multivariate failure time variables. In particular, we can test hypotheses for linear combinations of the  $\beta_k$ 's. The multivariate general linear hypothesis is written as

$$H_0: C\beta_T = 0,$$

where the  $r \times pK$  matrix  $C$  is called the contrast matrix. For example, if we want to test the hypothesis that the multivariate failure times do not depend on any covariates, then  $C$  will be the  $pK \times pK$  identity matrix. The Wald statistic for testing  $H_0$  is

$$(C\hat{\beta}_T)'(C\hat{Q}C')^{-1}(C\hat{\beta}_T),$$

which has an asymptotic  $\chi^2$  distribution with  $r$  degrees of freedom.

Next, suppose that we are interested in the effects of a particular covariate on the  $K$  failure time variables. Let us denote these  $K$  parameters by  $\eta_k$  ( $k = 1, \dots, K$ ). The  $\eta_k$ 's are obtained from  $\beta_T$  through a contrast matrix  $C$  for which  $C\beta_T = (\eta_1, \dots, \eta_K)'$ . If we assume that  $\eta_1 = \dots = \eta_K = \eta$ , it is natural to estimate  $\eta$  by a linear combination of the  $\hat{\eta}_k$ 's, that is,  $\sum_{k=1}^K h_k \hat{\eta}_k$  with  $\sum_{k=1}^K h_k = 1$ . The estimator  $\hat{\eta}$  with the array of weights  $(h_1, \dots, h_K)' = (e' \hat{\Psi}^{-1} e)^{-1} \hat{\Psi}^{-1} e$ , where  $e = (1, \dots, 1)'$  and  $\hat{\Psi} = C\hat{Q}C'$ , has the smallest asymptotic variance among all the linear estimators. It is obvious that the variance of  $\hat{\eta}$  can be estimated by  $(e' \hat{\Psi}^{-1} e)^{-1}$ . In applications, even if the  $\eta_k$ 's are unequal, we may still combine the  $\hat{\eta}_k$ 's to draw a conclusion about the 'average effect' of the covariate provided that there are no qualitative differences among the  $\eta_k$ 's.

### 3. Computer program

#### 3.1. General description

The MULCOX computer program was written in standard FORTRAN-77 with double arithmetic precision. The source program consists of 980

lines of codes and requires 29 kbytes of disk storage. No external subroutines or functions are used. The program can run on a mainframe computer or on a microcomputer.

The amount of CPU time used by MULCOX depends on the computer installation and the size of data. In general, the time consumption is minimal on a mainframe even for large data sets.

The program allows arbitrary values of  $n$ ,  $p$  and  $K$ . The matrices of data and computational results are stored in a single one-dimensional array  $A$ . The dimension of  $A$  may be modified by the user if necessary.

The covariates to be included in the model can be different from the prognostic variables in the data file. There is a subroutine in MULCOX which can be easily modified for necessary data transformation. The user who is unfamiliar with FORTRAN programming should transform the data through a software of his/her choice before running MULCOX.

#### 3.2. Input

The program MULCOX requires two separate groups of input: the data input and the control parameters input. The data (times, failure indicators and covariates) should be in the form of Table 1. Note that we have deliberately labeled the covariates in Table 1 by the  $Z_{ki}^*$ 's to distinguish them from the covariates to be included in the model (the  $Z_{ki}$ 's). Note also that the number of covariates in the data file, say,  $q$  can be different from the number of covariates in the model  $p$ .

The control parameters are described in Table 2. These parameters can be read from the keyboard upon execution of MULCOX or from an input file.

When specifying the format of the data file, the user may enter FREE or free if the data items are separated by spaces or commas; otherwise, a FORTRAN format expression with real and skip fields such as (5X, F10.5, F5.1, F8.5, F5.1, 3X, 2F6.3) is required.

#### 3.3. Output

The computational results are written to the output file specified by the user. The output consists

TABLE 1

The structure of data input

$X_{11}$	$\Delta_{11}$	$X_{21}$	$\Delta_{21}$	...	$X_{K1}$	$\Delta_{K1}$	$Z_{11}^*$	$Z_{21}^*$	...	$Z_{q1}^*$
$X_{12}$	$\Delta_{12}$	$X_{22}$	$\Delta_{22}$	...	$X_{K2}$	$\Delta_{K2}$	$Z_{12}^*$	$Z_{22}^*$	...	$Z_{q2}^*$
$\vdots$					$\vdots$					$\vdots$
$X_{1n}$	$\Delta_{1n}$	$X_{2n}$	$\Delta_{2n}$	...	$X_{Kn}$	$\Delta_{Kn}$	$Z_{1n}^*$	$Z_{2n}^*$	...	$Z_{qn}^*$

TABLE 2

The input of control parameters

Parameter	Type
title	character
file name of data input	character
file name of program output	character
$n$	integer
$K$	integer
name of the 1st failure time variable	character
$\vdots$	$\vdots$
name of the $K$ th failure time variable	character
$q$	integer
$p$	integer
name of the 1st covariate	character
$\vdots$	$\vdots$
name of the $p$ th covariate	character
format of data input	character
number of multivariate hypotheses	integer
row dimension of $C$ for the first hypothesis	integer
$\vdots$	$\vdots$
row dimension of $C$ for the last hypothesis	integer
number of common parameters to be estimated	integer
row dimension of $C$ for the first common parameter	integer
$\vdots$	$\vdots$
row dimension of $C$ for the last common parameter	integer
matrix $C$ for the first multivariate hypothesis	real
$\vdots$	$\vdots$
matrix $C$ for the last multivariate hypothesis	real
matrix $C$ for the first common parameter	real
$\vdots$	$\vdots$
matrix $C$ for the last common parameter	real

of four parts: I. estimation of marginal models, II. estimation of joint covariance matrix, III. testing multivariate hypotheses, and IV. estimation of common parameters. The output is self-explanatory.

#### 4. Application

To illustrate the use of MULCOX, let us consider a recent clinical trial evaluating the effectiveness of the drug ribavirin for treating patients with acquired immunodeficiency syndrome (AIDS). Thirty-six patients were randomly assigned to one of three groups: placebo, low-dose ribavirin and high-dose ribavirin. One of the main objectives of the study was to investigate the antiretroviral capability of ribavirin over time. Serum samples of each patient were collected at weeks 4, 8 and 12. The HIV-1 virus expression was evaluated by recording the number of days a patient's lymphocytes were in culture before virus positivity was detected. Hence, each patient should have three such event times. Some observations were missing, however, because patients did not make the scheduled visits or because serum specimens were inadequate for laboratorial analysis. In addition, censored observations occurred when the culture required a longer period of time to register as virus positive than was achievable in the laboratory, or when the serum sample was contaminated before positivity was detected.

In this example,  $V_{ki}$  is the number of days to virus positivity in the  $k$ th serum sample of the  $i$ th patient ( $k = 1, 2, 3; i = 1, \dots, 36$ ). Let  $Z_{1i} = 1$  if the  $i$ th patient was in the low-dose group and  $Z_{1i} = 0$  otherwise, and let  $Z_{2i} = 1$  if the  $i$ th patient was in the high-dose group and  $Z_{2i} = 0$  otherwise. The corresponding regression coefficients  $\beta_{1k}$  and  $\beta_{2k}$  can be interpreted, respectively, as the treatment effects from the low-dose and high-dose ribavirin after  $k$  months of treatment.

The data from this study are shown in Fig. 1. The last column is the treatment label: '1', '2' and '3' denote placebo, low-dose ribavirin and high-dose ribavirin, respectively. The subroutine given in Fig. 2 was used to create covariates  $Z_{1i}$  and  $Z_{2i}$  from the treatment label.

9	1	6	1	7	1	1
6	1	4	1	5	1	2
21	1	9	1	8	0	3
13	1	7	1	21	0	3
31	1	19	0	21	0	2
16	1	6	1	20	1	3
16	1	17	1	21	0	2
4	1	5	1	10	1	1
6	1	7	1	6	1	1
3	1	8	1	6	1	3
10	1	0	0	21	0	1
27	0	19	0	0	0	2
7	1	16	1	23	0	2
21	1	0	0	25	0	3
15	1	8	1	0	0	1
3	1	0	0	6	1	1
28	0	7	1	19	0	2
7	1	19	1	3	1	3
28	0	3	1	16	1	2
4	1	7	1	3	1	1
15	1	12	1	16	1	2
11	1	13	1	21	0	3
27	0	18	0	9	1	3
14	1	14	1	6	1	3
8	1	11	1	15	1	3
18	1	21	0	22	1	2
9	1	12	1	12	1	1
8	1	4	1	7	1	3
9	1	19	1	19	0	1
8	1	3	1	9	1	3
6	1	5	1	6	1	1
9	1	0	0	18	1	1
8	1	4	1	7	1	2
9	1	20	0	17	0	1
19	0	10	1	17	0	3
4	1	21	0	7	1	2

Fig. 1. Data input of the sample run: `aids.dat`.

The control parameters for this run were provided directly from a computer terminal (see Fig. 3). The two multivariate hypotheses to be tested were  $H_0: \beta_{11} = \beta_{12} = \beta_{13}$  and  $H_0: \beta_{21} = \beta_{22} = \beta_{23}$ . The one common parameter to be estimated was  $\eta = \beta_{11} = \beta_{12} = \beta_{13}$ .

This run only took a couple of seconds on a VAX-8550 computer. Its output is displayed in Fig. 4. The results indicated that high-dose ribavirin was beneficial to AIDS patients only at week 4. Although the effects of low-dose ribavirin also seemed to diminish over time, the observed



```

[jimmy]% MULCOX
PLEASE ENTER THE TITLE OF THIS RUN
the study of ribavirin on AIDS patients
ENTER THE NAME OF DATA FILE
aids.dat
ENTER THE NAME OF OUTPUT FILE
aids.out
ENTER THE NUMBER OF STUDY SUBJECTS
36
ENTER THE NUMBER OF FAILURE TIME VARIABLES
3
ENTER THE NAME OF FAILURE TIME VARIABLE 1
Week 4
ENTER THE NAME OF FAILURE TIME VARIABLE 2
Week 8
ENTER THE NAME OF FAILURE TIME VARIABLE 3
Week 12
ENTER THE NUMBER OF COVARIATES IN THE DATA FILE
1
ENTER THE NUMBER OF COVARIATES IN THE ASSUMED MODEL
2
ENTER THE NAME OF COVARIATE 1
Low Dose
ENTER THE NAME OF COVARIATE 2
High Dose
ENTER THE FORMAT OF THE DATA FILE
free
ENTER THE NUMBER OF MULTIVARIATE HYPOTHESES
2
ENTER THE DIMENSION OF MULTIVARIATE HYPOTHESIS 1
2
ENTER THE DIMENSION OF MULTIVARIATE HYPOTHESIS 2
2
ENTER THE NUMBER OF COMMON PARAMETERS
1
ENTER THE NUMBER OF COMPONENTS IN COMMON PARAMETER 1
3
ENTER THE CONTRAST MATRIX FOR MULTIVARIATE HYPOTHESIS 1
1 0 -1 0 0 0
1 0 0 0 -1 0
ENTER THE CONTRAST MATRIX FOR MULTIVARIATE HYPOTHESIS 2
0 1 0 -1 0 0
0 1 0 0 0 -1
ENTER THE CONTRAST MATRIX FOR COMMON PARAMETER 1
1 0 0 0 0 0
0 0 1 0 0 0
0 0 0 0 1 0
PLEASE WAIT !
[jimmy]%

```

Fig. 3. Control parameters input of the sample run.

```

* * * * *
*
*   Regression Analysis with Multiple Failure Time
*
*   Variables Based on Cox Proportional Hazards Models
*
* * * * *

```

- References: 1. L.J. Wei, D. Y. Lin and L. Weissfeld (1989).  
 Regression analysis of multivariate incomplete  
 failure time data by modeling marginal distributions.  
 Journal of the American Statistical Association  
 84, 1065-1073.
2. D. Y. Lin (1990). MULCOX: a computer program for  
 the Cox regression analysis of multiple failure  
 time variables. Computer Methods and Programs in  
 Biomedicine (in press).

PROBLEM TITLE IS: the study of ribavirin on AIDS patients  
 DATA FILE IS: aids.dat  
 OUTPUT FILE IS: aids.out

I. ESTIMATION OF MARGINAL MODELS  
 -----

FAILURE TIME VARIABLE = Week 4

Total number of study subjects = 36  
 Number of missing observations = 0  
 Number of observed failure times = 31

Log partial likelihood with zero beta = -90.17  
 Maximum log partial likelihood = -86.29  
 Global chi-square tests with D.F. = 2  
 -2 log L.R. = 7.75 P-value = 0.02079  
 Score = 8.47 P-value = 0.01450

PARAMETER	ESTIMATE	STANDARD ERROR	EST/S.E.
-----	-----	-----	-----
Low dose	-1.39393	0.52493	-2.65544
High dose	-0.93831	0.45518	-2.06140

Fig. 4. Program output of the sample run: aids.out.



FAILURE TIME VARIABLE = Week 8

Total number of study subjects = 36  
 Number of missing observations = 4  
 Number of observed failure times = 26

Log partial likelihood with zero beta = -75.46  
 Maximum log partial likelihood = -74.27  
 Global chi-square tests with D.F. = 2  
 -2 log L.R. = 2.39 P-value = 0.30264  
 Score = 2.27 P-value = 0.32098

PARAMETER	ESTIMATE	STANDARD ERROR	EST/S.E.
Low dose	-0.65523	0.52309	-1.25261
High dose	0.01999	0.46868	0.04266

FAILURE TIME VARIABLE = Week 12

Total number of study subjects = 36  
 Number of missing observations = 2  
 Number of observed failure times = 22

Log partial likelihood with zero beta = -66.80  
 Maximum log partial likelihood = -66.18  
 Global chi-square tests with D.F. = 2  
 -2 log L.R. = 1.26 P-value = 0.53380  
 Score = 1.28 P-value = 0.52756

PARAMETER	ESTIMATE	STANDARD ERROR	EST/S.E.
Low dose	-0.61512	0.55440	-1.10953
High dose	-0.33102	0.50498	-0.65551

II. ESTIMATION OF JOINT COVARIANCE MATRIX

0.245e+00	0.753e-01	0.507e-01	0.168e-01	0.107e+00	0.607e-01
0.753e-01	0.136e+00	0.265e-01	0.406e-01	0.458e-01	0.911e-01
0.507e-01	0.265e-01	0.287e+00	0.119e+00	0.133e+00	0.907e-01
0.168e-01	0.406e-01	0.119e+00	0.167e+00	0.763e-01	0.686e-01
0.107e+00	0.458e-01	0.133e+00	0.763e-01	0.257e+00	0.114e+00
0.607e-01	0.911e-01	0.907e-01	0.686e-01	0.114e+00	0.229e+00

Fig. 4 (continued).

### III. TESTING MULTIVARIATE HYPOTHESES

---

#### MULTIVARIATE HYPOTHESIS 1

The contrast matrix is as follows:

1.	0.	-1.	0.	0.	0.
1.	0.	0.	0.	-1.	0.

Wald statistic = 2.18402  
 Degrees of freedom = 2  
 P-value = 0.33554

#### MULTIVARIATE HYPOTHESIS 2

The contrast matrix is as follows:

0.	1.	0.	-1.	0.	0.
0.	1.	0.	0.	0.	-1.

Wald statistic = 4.67879  
 Degrees of freedom = 2  
 P-value = 0.09639

### IV. ESTIMATION OF COMMON PARAMETERS

---

#### COMMON PARAMETER 1

The contrast matrix is as follows:

1.	0.	0.	0.	0.	0.
0.	0.	1.	0.	0.	0.
0.	0.	0.	0.	1.	0.

The array of optimal weights is as follows:

0.44109    0.33964    0.21927

Estimator = -0.97227  
 Standard error = 0.38595  
 Z-score = -2.51917  
 Two-sided p-value = 0.01176

Fig. 4 (continued).

grateful to a referee, and Professors Morton B. Brown and L.J. Wei for their useful suggestions.

### References

- [1] N. Breslow, Covariance analysis of censored survival data, *Biometrics* 30 (1974) 89–99.
- [2] D.R. Cox, Regression models and life tables (with discussion), *J. R. Stat. Soc. B* 34 (1972) 187–220.
- [3] D.Y. Lin and L.J. Wei, The robust inference for the Cox proportional hazards model, *J. Am. Stat. Assoc.* 84 (1989) 1074–1078.
- [4] L.J. Wei, D.Y. Lin and L. Weissfeld, Regression analysis of multivariate incomplete failure time data by modeling marginal distributions, *J. Am. Stat. Assoc.* 84 (1989) 1065–1073.