

### B.3 Survival regression models and two-sample testing

1. (a) In a Weibull model, the survival function is  $S(x) = e^{-(\rho x)^\alpha}$ . Thus  $\log(-\log S(x)) = \alpha \log \rho + \alpha \log x$ , and if we plot  $\log(-\log \hat{S}(x))$  against  $\log x$  we should see something close to a straight line. Since the exponential model is a submodel of the Weibull (with  $\alpha = 1$ ), we can apply the likelihood ratio test. If  $\ell(\rho, \alpha)$  is the log likelihood, we have under the null model (that the data were sampled from an exponential distribution)

$$\sup_{(\rho, \alpha)} \ell(\rho, \alpha) - \sup_{\rho} \ell(\rho, 1) \sim \chi_1^2.$$

For the log-logistic model, we expect the plot of  $\log(\frac{1}{\hat{S}(x)} - 1)$  against  $\log x$  to be approximately linear.

- (b) Let  $S_1$  and  $S_2$  be the survival curves for the two populations, and  $S_0$  the baseline survival. Under the accelerated lifetime model,  $S_i(x) = S_0(\rho_i x)$  for some positive constants  $\rho_1, \rho_2$ . Then if we plot  $S_i(x)$  against  $\log x$ , we see that whatever value  $S_0$  takes at ordinate  $\log x$ ,  $S_i$  will take the same value at an interval of  $\log \rho_i$ . (The same will be true of any function of  $S_i$ .) Thus, the graphs corresponding to  $\hat{S}_1$  and  $\hat{S}_2$  should differ approximately by a uniform horizontal shift.

The proportional hazards assumption is best tested by plotting  $\log(-\log \hat{S}_i(x))$ . Under PH,  $S_i(x) = S_0(x)^{\rho_i}$ , which implies that

$$\log(-\log S_i(x)) = \log(-\rho_i \log S_0(x)) = \log(\rho_i) + \log(-\log S_0(x)).$$

Thus, if  $\log(-\log \hat{S}_i(x))$  is plotted against  $x$ , the two graphs should differ approximately by a constant vertical shift if the two groups satisfy the PH assumption. The same is true if we plot  $\log(-\log \hat{S}_i(x))$  against any function of  $x$ . Thus, if we plot  $\log(-\log \hat{S}_i(x))$  against  $\log x$ , we will see a constant vertical shift reflecting the PH assumption, and a constant horizontal shift reflecting the AL assumption.

- (c) The computations for the Kaplan–Meier estimator are given in Table [B.1](#). In figure [B.1](#) we plot the two survival curves (red for control, black for treatment), as  $\log(-\log \hat{S})$  against  $\log x$ . Both look reasonably close to lines, so it would be reasonable to suppose that they came from Weibull models. The lines are approximately parallel, suggesting that the  $\alpha$  parameters are approximately the same. This means that one curve may be obtained from another by a horizontal or vertical shift, suggesting that PH or AL would be appropriate. (Weibull curves with the same  $\alpha$  parameter, it should be noted, satisfy both hypotheses.)

$t_j$	$d_j$	$n_j$	$\hat{h}_j$	$\hat{S}(t_j)$
1	2	21	0.095	0.905
2	2	19	0.105	0.810
3	1	17	0.059	0.762
4	2	16	0.125	0.667
5	2	14	0.143	0.572
8	4	12	0.333	0.381
11	2	8	0.250	0.286
12	2	6	0.333	0.191
15	1	4	0.250	0.143
17	1	3	0.333	0.095
22	1	2	0.500	0.048
23	1	1	1.000	0.000

$t_j$	$d_j$	$n_j$	$\hat{h}_j$	$\hat{S}(t_j)$
6	3	21	0.143	0.857
7	1	17	0.059	0.806
10	1	15	0.067	0.752
13	1	12	0.083	0.690
16	1	11	0.091	0.627
22	1	7	0.143	0.537
23	1	6	0.167	0.448

Table B.1: Estimates for control group (left) and treatment group (right) in Gehan study.

We test the hypothesis by finding maximum likelihood estimators. The log likelihood for the exponential distribution are

$$\ell(\lambda) = \sum_i (-\lambda x_i) + d \log \lambda,$$

where  $d$  is the number of uncensored observations. Since the maximum likelihood estimator is  $\hat{\lambda} = d / \sum x_i$ , we get maximum likelihoods of

$$\ell_{exp}^* = d \left( \log d - 1 - \log \sum x_i \right).$$

For the Weibull distribution we have

$$\ell(\rho, \alpha) = - \sum (\rho x_i)^\alpha + d(\alpha \log \rho + \log \alpha) + \sum_{i \text{ uncensored}} (\alpha - 1) \log x_i.$$

There is no closed form solution, but we can optimise numerically, yielding estimates

	Treatment	Control
$\hat{\lambda}$	0.025	0.12
$\ell_{exp}^*$	-42.17	-66.35
$\hat{\rho}$	0.030	0.11
$\hat{\alpha}$	1.35	1.37
$\ell_{weib}^*$	-41.66	-64.92

The log likelihood ratio for the treatment group is thus  $(-41.66) - (-42.17) = 0.51$ , and for the control group it is 1.43. Comparing these to the  $\chi^2$  distribution with 1 degree of freedom, we see that the cutoff for rejecting the null hypothesis that  $\alpha = 1$  at the 0.05 significance level would be 3.84. Thus, we cannot reject the null hypothesis for either group.

2. (a) Assuming no ties, the partial likelihood is constructed by computing the probability that the subjects failed in exactly the order observed, conditioned on the times observed.

The proportional hazards (PH) assumption says that subject  $i$  has hazard rate  $h_i(x) = r_i h_0(x)$  at time  $x$ , where  $h_0$  is an unspecified baseline hazard. In the regression approach, we think of  $r_i$  as a function  $r(y_i)$  of a vector  $y_i$  of covariates. The linear approach is to suppose  $\phi(r(y)) = \beta \cdot y$ , where  $\phi$  is the *link function* and  $\beta$  is a vector of parameters to estimate. In the Cox model we use the logarithmic link function, so that  $r(y) = e^{\beta \cdot y}$ . The partial likelihood is defined as

$$L_P(\beta; y) := \prod_{t_i} \frac{e^{\beta y_{(i)}}}{\sum_{j \in R_j} e^{\beta y_j}},$$

where  $x_{(i)}$  represents the covariates of the subject failing at time  $t_i$  and  $R_i$  is the *risk set*, of those subjects at risk at  $t_i$ .

We use  $L_P$  as though it were a likelihood. We compute the parameters  $\hat{\beta}$  that maximise  $L_P$ . Under the assumption that the observations came from the distribution given by this model with some (unknown) parameter  $\beta$ , the estimate  $\hat{\beta}$  is asymptotically normal, with mean  $\beta$  and variance matrix that may be estimated by

$$\left[ E \left( - \frac{\partial^2 \ell_P}{\partial \beta \partial \beta^T} \right) \right]^{-1}, \text{ where } \ell_P = \log L_P.$$

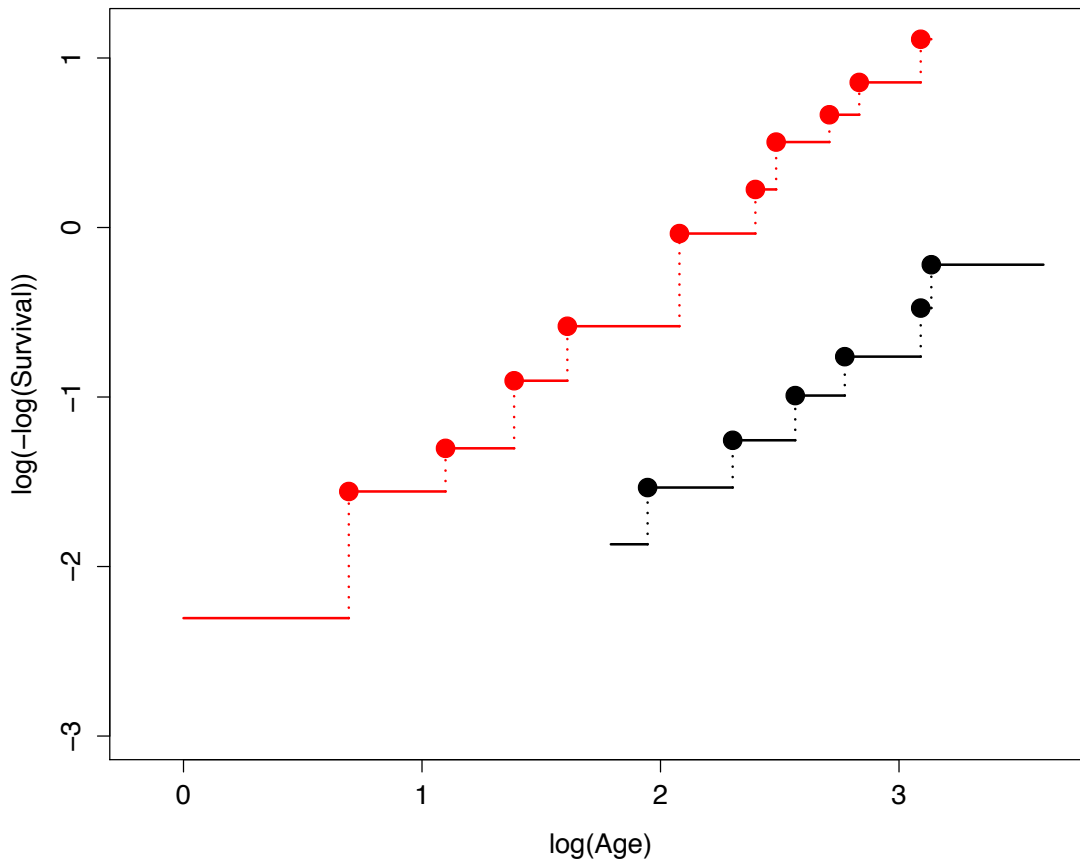


Figure B.1: Plot of estimated survival for Gehan leukaemia data. The control group is in red, the treatment group is black.

(b) The hazard ratio is

$$\begin{aligned} \frac{h(\text{clinic} = 1, \text{prison}=0)}{h(\text{clinic} = 0, \text{prison}=1)} &= \frac{e^{\hat{\beta} \cdot y_1}}{e^{\hat{\beta} \cdot y_2}} \\ &= \frac{e^{-1.009}}{e^{0.327}} \\ &= 0.263. \end{aligned}$$

(c) The log hazard ratio for prison/no prison is 0.327, with standard error 0.167. A 95% confidence interval for the coefficient is  $0.327 \pm 1.96 \cdot 0.167 = (0.0, 0.654)$ . Thus a 95% confidence interval for the hazard ration is  $e^{(0.0, 0.654)} = (1.00, 1.92)$ .

- We give below R code for computing this in two different ways: Using the function `survdif`, which does the computation automatically, and by extracting the relevant quantities from the survival object and doing the computation directly.

We get  $Z = -1.67$ , which corresponds to a  $p$ -value of 0.09.

Using `survdif` we get the same result, but it is reported as a chi-squared statistic of 2.8 (which is  $1.67^2$ ) on 1 degree of freedom.

## SURVDIFF CODE

```

> require('survival')
> require('KMsurv')
> data(tongue)
> attach(tongue)

>
> tongue.surv=Surv(time,delta)
> tongue.fit=survfit(tongue.surv~type)
> tdiff=survdiff(tongue.surv~type)
> tdiff
Call:
survdiff(formula = tongue.surv ~ type)

N Observed Expected (O-E)^2/E (O-E)^2/V
type=1 52      31      36.6      0.843      2.79
type=2 28      22      16.4      1.873      2.79

Chisq= 2.8  on 1 degrees of freedom, p= 0.0949

```

## DIRECT COMPUTATION

```

# Problem sheet 4, question 1
require('survival')
require('KMsurv')
data(tongue)
attach(tongue)

tongue.surv=Surv(time,delta)
tongue.fit=survfit(tongue.surv~type)

n1=tongue.fit$strata[1]
n2=tongue.fit$strata[2]

# Input two vectors of times t1,t2, and
# numbers at risk n1,n2 whose length is 1 longer than the t's
# Output four vectors I1, I2, (of same length as t1,t2) and Y1,Y2
# I1[k] gives an index of I2 corresponding to
# the last time in t2 that precedes t1[k]
# Thus, we have t2[I1[k]]<=t1[k] < t2[I1[k]+1],
# and r2[I1[k]+1] is the number of type 2 individuals at risk
# at the time t1[k] (when there are r1[k] type 1 individuals)
# Y1=r1[I1]

crossrisk=function(t1,t2,r1,r2){
  I1=rep(0,length(t1))
  I2=rep(0,length(t2))
  for(i in seq(length(t1))){
    I1[i]=1+sum(t1[i]>t2)
  }
  for(i in seq(length(t2))){

```

```

    I2[i]=1+sum(t2[i]>t1)
  }
  list(I1,I2,r1[I2],r2[I1])
}

r1=tongue.fit$n.risk[seq(n1)]
r2=tongue.fit$n.risk[seq(n1+1,n1+n2)]

r1=c(r1,r1[n1]-tongue.fit$n.event[n1]-tongue.fit$n.censor[n1])
r2=c(r2,r2[n2]-tongue.fit$n.event[n1+n2]-tongue.fit$n.censor[n1+n2])
t1=tongue.fit$time[seq(n1)]
t2=tongue.fit$time[seq(n1+1,n1+n2)]

cr=crossrisk(t1,t2,r1,r2)

Y1=c(r1[-n1],cr[[3]])
Y2=c(cr[[4]],r2[-n2])
# Note: r1 and r2 had an extra count added on to make crossrisk work
d1=c(tongue.fit$n.event[seq(n1)],rep(0,n2))
d2=c(rep(0,n1),tongue.fit$n.event[seq(n1+1,n1+n2)])

t=c(t1,t2)

# We have to deal with the problem of ties between times for the two groups

dup1=which(duplicated(t,fromLast=TRUE))
dup2=which(duplicated(t))
ndup=length(dup1)

# Type 2 Event counts are removed from the second appearance
# and placed in the first appearance
d2[dup1]=d2[dup2]
d2=d2[-dup2]
d1=d1[-dup2]

# Type 2 at-risk counts are removed from the second appearance
# and placed in the first appearance
Y2[dup1]=Y2[dup2]
Y2=Y2[-dup2]
Y1=Y1[-dup2]
t=t[-dup2]

tord=order(t)
t=t[tord] #put times in order
## Now put everything else in the same order
Y=Y[tord]
Y1=Y1[tord]
Y2=Y2[tord]
d=d[tord]
d1=d1[tord]
d2=d2[tord]

Y=Y1+Y2

```

```

d=d1+d2

# Product of number at risk
atriskprod=Y1*Y2
includes=(atriskprod>0)&(d>0)
# We only get contributions if someone's at risk and events occurred at that time

Y=Y[includes]
Y1=Y1[includes]
Y2=Y2[includes]
d=d[includes]
d2=d2[includes]
d1=d1[includes]

t=t[includes]

wLR=Y1*Y2/Y
p=1
q=0

S=c(1,cumprod((Y-d)/Y))[-length(Y)] #K-M estimator for survival
wFH=(1-S)^q*S^p*wLR

# Now compute the test statistic

w=wLR

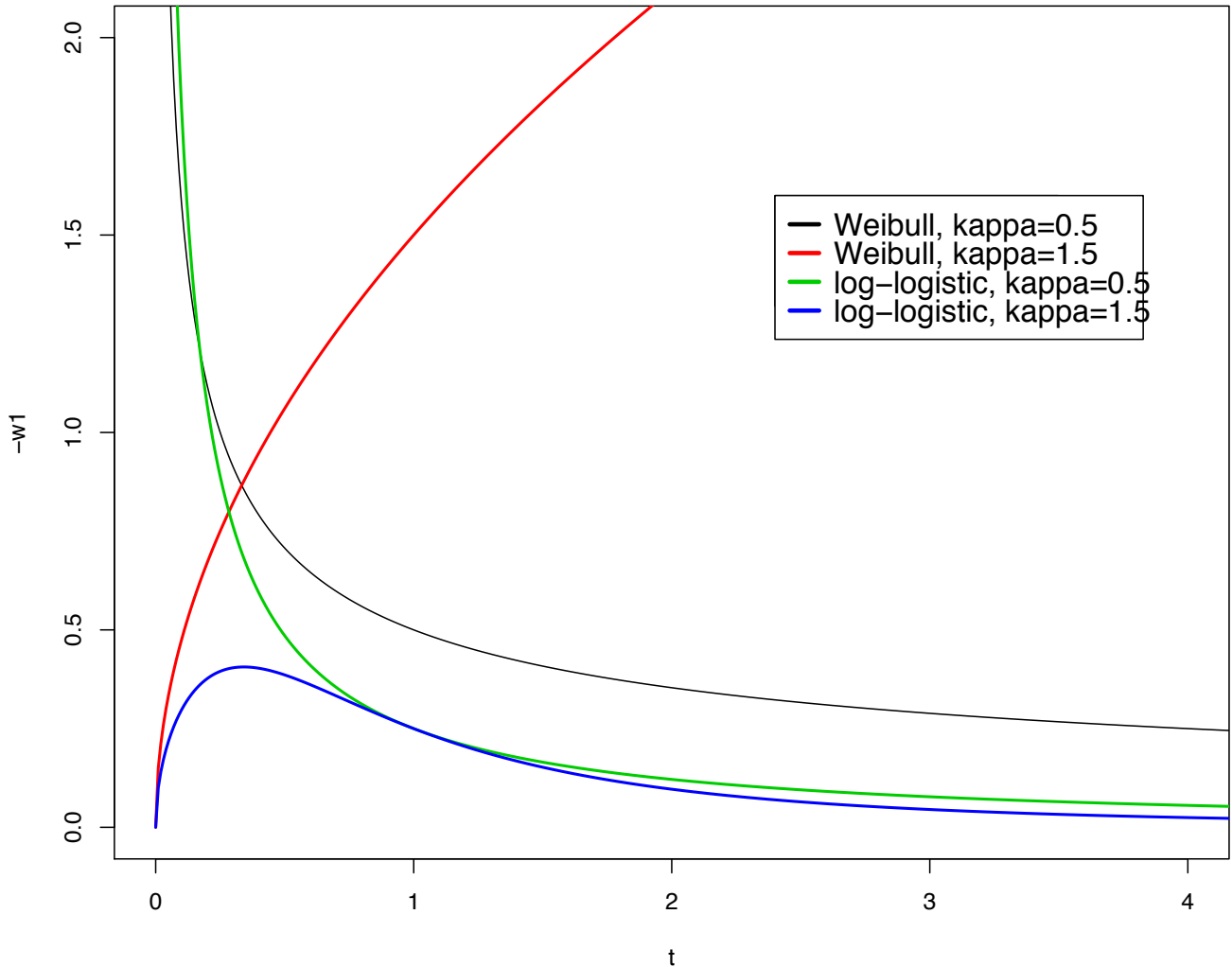
M=w*(d1/Y1-d2/Y2)
sigma=w*w*d*(Y-d)/Y2/Y1/(Y-1)
sK=d*Y1*Y2*(Y-d)/Y^2/(Y-1)

Z=sum(M)/sqrt(sum(sigma))

> Z
[1] -1.670246

```

4. (a) The plot is:



(b) One could consider using the log-logistic or the log-normal.

(c) The hazard function is  $h(t) = \kappa(\rho e^{\beta \cdot x})^\kappa t^{\kappa-1}$  and the survival function is  $e^{-(\rho e^{\beta \cdot x} t)^\kappa}$ . Hence the log likelihood is

$$\ell(\rho, \kappa, \beta) = n \log \kappa + n \kappa \log \rho + \kappa \sum \beta \cdot x_i + (\kappa - 1) \sum \log t_i - \sum (\rho e^{\beta \cdot x_i} t_i)^\kappa.$$

The MLE must satisfy

$$\begin{aligned} 0 &= \frac{\partial \ell}{\partial \rho} = \frac{n \kappa}{\rho} - \kappa \rho^{\kappa-1} \sum (e^{\beta \cdot x_i} t_i)^\kappa, \\ 0 &= \frac{\partial \ell}{\partial \kappa} = n \left( \frac{1}{\kappa} + \log \rho \right) + \sum \beta \cdot x_i + \sum \log t_i - \kappa \sum (\rho e^{\beta \cdot x_i})^\kappa t_i^{\kappa-1}, \\ 0 &= \frac{\partial \ell}{\partial \beta_j} = \kappa \sum_i x_{ij} (1 - (\rho e^{\beta \cdot x_i} t_i)^\kappa). \end{aligned}$$

Asymptotically, the estimators will be normally distributed. If some observations are right-censored, the log likelihood becomes

$$\ell(\kappa, \beta) = n_d \log \kappa + n_d \kappa \log \rho + \kappa \sum \delta_i \beta \cdot x_i + (\kappa - 1) \sum \delta_i \log t_i - \sum \left( \rho e^{\beta \cdot x_i} t_i \right)^\kappa$$

where  $n_d$  is the number of (uncensored) events observed.

5. (a) The times  $t_i$  are 50, 52, 58, 61, 67, 68, 70, 72, 75. A full description of the risk sets requires that we describe exactly which individuals are at risk. We number the males as  $M1, \dots, M12$  and  $F1, \dots, F12$ . We have then the risk sets

$$R_1 = \{M1, M3, M4, M5, M6, M7, M8, M9, M10, M11, M12, F1, F2, F3, F4, F5, F6, F7, F8, F9, F10, F11, F12\}$$

$$R_2 = \{M3, M4, M5, M6, M7, M8, M9, M10, M11, M12, F1, F2, F3, F4, F5, F6, F7, F8, F9, F10, F11, F12\}$$

$$R_3 = \{M4, M5, M6, M7, M8, M9, M10, M11, M12, F1, F3, F4, F5, F6, F7, F8, F9, F10, F11, F12\}$$

$$R_4 = \{M4, M5, M6, M7, M8, M9, M10, M11, M12, F1, F3, F4, F5, F6, F7, F8, F9, F10, F12\}$$

$$R_5 = \{M4, M5, M6, M7, M8, M9, M10, M12, F3, F4, F5, F6, F7, F8, F9, F10, F12\}$$

$$R_6 = \{M4, M5, M6, M7, M9, M10, M12, F3, F4, F5, F6, F7, F8, F9, F10, F12\}$$

$$R_7 = \{M5, M7, M9, M10, M12, F3, F4, F5, F7, F8, F9, F10, F12\}$$

$$R_8 = \{M5, M9, M10, M12, F4, F5, F9, F10, F12\}$$

$$R_9 = \{M10, M12, F4, F10\}.$$

Note that there is some ambiguity in breaking ties. When an observation is censored at time  $t_i$  we must decide whether to treat the censoring as having occurred just after or just before  $t_i$ : that is, was the individual available to have been counted if they had died at time  $t_i$  or not? We have chosen the former: Thus, for instance,  $R_9$  is the set of individuals at risk at time 75, and it includes M12, who was censored at age 75. Either one is acceptable — though details of the study may suggest one or the other interpretation — but it should be specified. Since we are interested only in the binary covariate of gender, we need only consider the risk sets as counting the numbers of males and females, coded as  $R_i = (m_i, f_i)$ . We may then summarise them as

$$R_1 = (11, 12) \quad R_2 = (10, 12) \quad R_3 = (9, 11) \quad R_4 = (9, 10) \quad R_5 = (8, 9) \\ R_6 = (7, 9) \quad R_7 = (5, 8) \quad R_8 = (4, 5) \quad R_9 = (2, 2).$$

- (b) Using the notation as above, and setting the vector of covariates to be  $x = (1, 0, 0, 1, 1, 1, 0, 0, 0)$  — coding female as 0 and male as 1 — we have the partial likelihood being

$$L_P = \prod_{i=1}^9 \frac{e^{\beta x_i}}{f_i + e^{\beta} m_i} = e^{4\beta} \prod_{i=1}^9 \left( f_i + e^{\beta} m_i \right)^{-1}. \quad (\text{B.9})$$

A plot of this function is in Figure [B.2](#). The maximum likelihood is attained at  $\beta = -0.042$ .



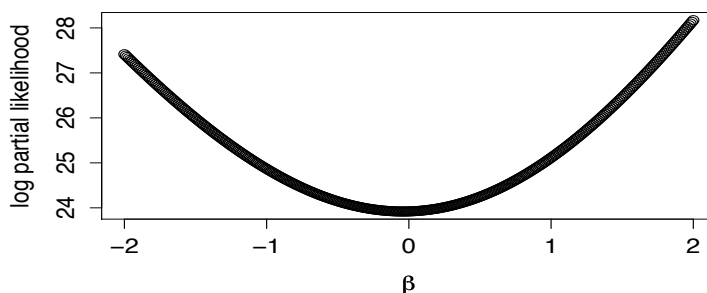


Figure B.2: Plot of negative logarithm of partial likelihood given by (B.9).

(c) We need to compute first  $\hat{S}$  for the combined population. We have

		event time								
		50	52	58	61	67	68	70	72	75
Male	$d_m$	1	0	0	1	1	1	0	0	0
	$m_i$	11	10	9	9	8	7	5	4	2
Female	$d_i^f$	0	1	1	0	0	0	1	1	1
	$f_i$	12	12	11	10	9	9	8	5	2
Total	$d_i$	1	1	1	1	1	1	1	1	1
	n	23	22	20	19	17	16	13	9	4
$\hat{S}(t_{i-1})$		1	0.957	0.913	0.867	0.822	0.773	0.725	0.669	0.595

Plugging these into the formula

$$Z = \frac{\sum_{i=1}^9 \left( d_i^m - n_i^m \frac{d_i}{n_i} \right)}{\sqrt{\sum_{i=1}^9 \frac{n_i^m n_i^f (n_i - d_i) d_i}{n_i^2 (n_i - 1)}}$$

we get  $Z = -.063$ , which should be like a draw from a normal distribution if the male and female survival times were drawn from the same distribution. In fact, we get a p-value of  $1 - 2\Phi(.063) = .95$ .

(d) For the Fleming–Harrington test we down-weight the later times, when very few are at risk, substituting

$$Z_{FH} = \frac{\sum_{i=1}^9 \hat{S}(t_{i-1}) \left( d_i^m - n_i^m \frac{d_i}{n_i} \right)}{\sqrt{\sum_{i=1}^9 \hat{S}(t_{i-1})^2 \frac{n_i^m n_i^f (n_i - d_i) d_i}{n_i^2 (n_i - 1)}}} = 0.105,$$

yielding a p-value for the two-sided test of 0.92. In either case, of course, we would not reject the null hypothesis. Of course, this is not surprising, as the sample is very small.

Note that this analysis could be improved by taking account of the pairing of twins.

(e) Death due to other causes is unlikely to be independent of CHD. Hence, non-informative censoring is questionable.

6. For clarity, we repeat the derivation in this particular setting. Let  $Z(t_j)$  be the  $2 \times 2$  matrix  $\mathbf{X}(t_j)^T \mathbf{X}(t_j)$ . Since  $x_i(t_j) = x_i = x_i^2$  is 0 or 1, we have

$$Z_{00} = \sum Y_i(t)^2 = n(t),$$

$$Z_{10} = Z_{01} = Z_{11} = \sum Y_i(t)x_i = n_1(t),$$

where  $n(t)$  is the number of individuals at risk at time  $t$ , and  $n_j(t)$  is the number of individuals at risk at time  $t$  with  $x_i = j$ . The determinant is  $n_1(t)n_0(t)$ . Inverting this, we get

$$\mathbf{X}^{-}(t) = \begin{pmatrix} \frac{1}{n_0(t)} & -\frac{1}{n_0(t)} \\ -\frac{1}{n_0(t)} & \frac{1}{n_1(t)} + \frac{1}{n_0(t)} \end{pmatrix} \begin{pmatrix} Y_1(t) & Y_2(t) & \cdots & Y_n(t) \\ x_1 Y_1(t) & x_2 Y_2(t) & \cdots & x_n Y_n(t) \end{pmatrix},$$

as long as  $n_1(t)n_0(t) > 0$  (and 0 otherwise). Thus, since  $Y_{i_j}(t_j)$  is always 1 (since the individual who has an event must, by definition, be at risk),

$$\mathbf{X}^{-}(t_j)_{cdotij} = \begin{pmatrix} \frac{1-x_0}{n_0(t_j)} \\ \left(-\frac{1-x_0}{n_0(t_j)} + \frac{x_0}{n_1(t_j)}\right) \end{pmatrix}$$

Hence the baseline cumulative hazard estimate is

$$\hat{B}_0(t) = \sum_{t_j \leq t} \frac{1-x_0}{n_0(t_j)} = \sum_{t_j \leq t: x_{i_j}=0} \frac{1}{n_0(t_j)},$$

which is the definition of the Nelson–Aalen estimator for the cumulative hazard, considering only the individuals with  $x_i = 0$ ; and the estimated cumulative increment due to  $x_i = 1$  is

$$\hat{B}_1(t) = \sum_{t_j \leq t: x_{i_j}=1} \frac{1}{n_1(t_j)} - \sum_{t_j \leq t: x_{i_j}=0} \frac{1}{n_0(t_j)},$$

7. (a) As described in the previous question, the difference may be estimated by the difference between the Nelson–Aalen estimators:

$$\hat{B}_1(t) = \hat{H}_0(t) - \hat{H}_1(t) = \sum_{t_j \leq t} \frac{d_{0j}}{n_{0j}} - \frac{d_{1j}}{n_{1j}}.$$

Calling the Maintenance group number 1, and Nonmaintenance number 0, we read off of Table 4.3  $\hat{H}_1(20) = \hat{H}_1(18) = 0.32$ , and  $\hat{H}_0(20) = 0.49$ , yielding

$$\hat{B}_1(20) = 0.17.$$

The variance will be the sum of the variances of the two estimators (since they are independent). As long as there are no ties between events from different groups, this may be estimated by

$$\sum_{t_j \leq t} \frac{d_{0j}}{n_{0j}^2} + \sum_{t_j \leq t} \frac{d_{1j}}{n_{1j}^2}.$$

From the table we can see that this is

$$\sigma_1^2(20) + \sigma_0^2(20) = \frac{1}{12^2} + \frac{1}{11^2} + \frac{1}{10^2} + \frac{1}{9^2} + \frac{1}{8^2} + \frac{1}{10^2} + \frac{1}{8^2} = 0.0788.$$

Thus, an approximate 95% confidence interval for  $\hat{B}_1(20)$  would be

$$0.17 \pm 0.28 \cdot 1.96 = 0.17 \pm 0.550.$$

(b) The Cox model fit by `coxph` produced the outcome

coxph(formula = Surv(time, status) ~ x, data = aml)					
	coef	exp(coef)	se(coef)	z	p
×Nonmaintained	0.916	2.5	0.512	1.79	0.074
Likelihood ratio test=3.38 on 1 df p=0.0658 n= 23					

In Table 5.2 we tabulated the estimators for the baseline hazard, obtaining  $\hat{H}_0(18) = 0.254$ . A central estimate for the difference in cumulative hazard between the two groups would be

$$(1 - e^{\hat{\beta}})\hat{A}_0(18) = 1.5 \cdot 0.254 = 0.38.$$

We see that this is a substantially larger estimate than we made in the nonparametric model. This is consistent with the plot in Figure 5.4, where the purple circles and blue crosses (representing the survival estimates from the proportional hazards model for the two groups) are further apart at  $t_j = 18$  than the black and red lines (representing the Kaplan–Meier estimators). This reflects that fact that the separate Kaplan–Meier estimators are cruder, making larger jumps at less frequent intervals.

To estimate the standard error, we begin by assuming (with little justification) that the estimators  $\hat{\beta}$  and  $\hat{H}_0(t_j)$  are approximately independent. Then we can use the delta method to estimate the variance. Let  $\sigma_{\hat{\beta}}^2$  be the variance of  $\hat{\beta}$ , and  $\sigma_H^2$  the variance of  $\hat{H}(18)$ . So we can represent

$$\hat{\beta} \approx \beta_0 + \sigma_{\beta}Z, \quad \hat{H}_0(18) = H_0(18) + \sigma_H Z',$$

where  $Z$  and  $Z'$  are standard normal (also approximately independent). We already have the estimate  $\hat{\sigma}_{\beta} \approx 0.512$ . We haven't given a formula for an estimator of  $\sigma_H(18)$ , but we can easily compute it with R.

```
require(survival)

cp=coxph(Surv(time,status)~x,data=aml)

aml.fit=survfit(cp)

aml.fit$std.err[aml.fit$time==18]
[1] 0.150247
```

Then our estimator for the difference in cumulative hazard is

$$\begin{aligned} (1 - e^{\hat{\beta}})\hat{H}_0(18) &\approx \left(1 - e^{\beta_0 + \sigma_{\beta}Z}\right) (H_0(18) + \sigma_H Z') \\ &\approx \left(1 - e^{\beta_0} (1 + \sigma_{\beta}Z)\right) (H_0(18) + \sigma_H Z') \\ &\approx \left(1 - e^{\beta_0}\right) H_0(18) - e^{\beta_0} \sigma_{\beta} H_0(18) Z + \left(1 - e^{\beta_0}\right) \sigma_H Z' - e^{\beta_0} \sigma_{\beta} \sigma_H Z Z'. \end{aligned}$$

(Note that the approximation in the first line is based on assuming  $\sigma_{\beta}$  is much smaller than  $\beta_0$ , which isn't really very true here.) As long as we are assuming independence of  $Z$  and  $Z'$ , the variance will be approximately

$$\left(e^{\beta_0} \sigma_{\beta} H_0(18)\right)^2 + \left(\left(1 - e^{\beta_0}\right) \sigma_H\right)^2 = 0.325^2 + .225^2 = 0.156,$$

so the standard error is about 0.395.

A better estimate, also taking into account the dependence between  $\hat{\beta}$  and  $\hat{H}_0$ , could be obtained by not using the delta method, but instead treating the normal distribution of  $\hat{\beta}$  as a Bayesian posterior distribution on  $\beta_0$ . For a range of possible  $\beta_0$  we can compute an approximate mean and variance for  $\hat{H}_0$ , and then compute a Monte Carlo estimator of the variance of  $\hat{\Gamma}$ .

- (c) We let  $x_0(t)$  be the covariate trajectory for this individual, so recalling that the maintained group is the baseline this means that

$$x_0(t) = \begin{cases} 0 & \text{if } t \leq 10, \\ 1 & \text{if } t > 10. \end{cases}$$

Using the formula (5.9) we estimate for this individual

$$\begin{aligned} \hat{H}(20 | x_0) &= \int_0^{20} e^{\beta x_0(u)} d\hat{H}_0(u) \\ &= \int_0^{10} d\hat{H}_0(u) + \int_{10}^{20} e^{\beta} d\hat{H}_0(u) \\ &\approx \hat{H}_0(10) + 2.5(\hat{H}_0(20) - \hat{H}_0(10)) \\ &= 0.14 + 2.5(0.114) \\ &= 0.425. \end{aligned}$$