# Problem Set 4

## Applied Stats II

## Due: April 16, 2023

## Instructions

- Please show your work! You may lose points by simply writing in the answer. If the problem requires you to execute commands in R, please include the code you used to get your answers. Please also include the .R file that contains your code. If you are not sure if work needs to be shown for a particular problem, please ask.

- Your homework should be submitted electronically on GitHub in .pdf form.

- This problem set is due before 23:59 on Sunday April 16, 2023. No late assignments will be accepted.

## Question 1

We're interested in modeling the historical causes of child mortality. We have data from 26855 children born in Skellefteå, Sweden from 1850 to 1884. Using the "child" dataset in the eha library, fit a Cox Proportional Hazard model using mother's age and infant's gender as covariates. Present and interpret the output.

## My Answer - Question 1

Looking at the variables in the imported data(infants)

- enter - Age (in days) of case when its mother died

- exit - Age (in days) at death or right censoring (at age 365 days)

- event - Follow up ends with death (1) or right censoring (0)

- mother - dead for cases, alive for controls

- age - Mother's age at infant's birth

- sex - The infant's sex

- parish - Birth parish, either Nedertornea or not Nedertornea (other)

- civst - Civil status of mother, married or unmarried

- ses - Socio economic status of mothe, either farmer or not farmer (other)

- year - Year of birth of the infant

In order to fit a Cox Proportional Hazard Model the coxph() function from the 'survival' package is required. This needs to be in a Surv() format which represents time to event data. The first argument is time variable.

Testing with Weibull proportional hazard model. Proportional hazards model with parametric baseline hazard(s). Allows for stratification with different scale and shape in each stratum, and left truncated and right censored data. Create a new variable called "gender"

in the "infants" data frame, and assign it values based on the "sex" variable. Specifically, the first command sets all the values of "gender" to missing (NA) for all observations in the "infants" data frame. The second command sets the value of "gender" to 0 for all observa-

tions in the "infants" data frame where the value of "sex" is equal to 'boy'. This means that all male infants in the data frame will have a value of 0 in the "gender" variable. The third

command sets the value of "gender" to 1 for all observations in the "infants" data frame where the value of "sex" is equal to 'girl'. This means that all female infants in the data frame will have a value of 1 in the "gender" variable. In summary, these commands create

a new variable "gender" in the "infants" data frame and assign values of 0 and 1 to it based on the "sex" variable. This is a common technique used to recode categorical variables with two categories into a binary format for statistical analysis. Kaplan Meier survival curve for

the "infants s" variable in the "infants" data frame and provide a summary of the survival probabilities at specific time points. This is a common technique used in survival analysis to estimate the probability of an event (e.g., death, failure) over time, based on a set of predictor variables.

```
1
2
3 # Read in the data and have a look at it
4 data(infants)
5 glimpse(infants)
6 infants
7 table(infants)
8
9 # inspecting the data
10 str(infants)
11 names(infants)
12
```

```
13  table(infants$sex)
14  # girl   boy
15  #  34    71
16
17  prop.table(table(infants$sex))
18  # girl        boy
19  # 0.3238095 0.6761905
20
21  # recode
22  infants$gender <- NA
23  infants$gender[infants$sex=='boy'] <- 0
24  infants$gender[infants$sex=='girl'] <- 1
25
26  # histogram to look at the visual of the age and sex variables
27  plot_ly(infants, x = ~age, color = ~sex) %>%
28    add_histogram()
29
30
31  fit <- coxreg(Surv(enter, exit, event) ~ strata(stratum) + mother, data
32              = infants)
33  fit
34  fit.w <- phreg(Surv(enter, exit, event) ~ mother + parish + ses, data =
35                  infants)
36  summary(fit.w) ## Weibull proportional hazards model.
37  plot(fit.w)
38
39  coef(fit.w) # for extracting the coefficients from the model
40
41
42  infants_s <- with(infants, Surv(enter, exit, event))
```

The Kaplan-Meier method is a non parametric statistic that allows you to estimate the survival function.

```
1
2
3  # The Kaplan-Meier method is a non parametric statistic that allows you to
       estimate the survival function
4
5  # plotting Kaplan-Meier method
6  # The resulting object "kaplan" is a "survfit" object that stores the Kaplan-
       Meier estimates of the survival probabilities,
7  # as well as the number of observations and the number of events at each time
       point.
8  # The object can be used for further analysis or plotting.
9
10
11  kaplan <- survfit(infants_s ~ 1, data = infants)
12  summary(kaplan, times = seq(0, 15, 1))
13
14  plot(kaplan, main = "Kaplan-Meier Method", xlab = "Days", ylim = c(0.5, 1))
15  autoplot(kaplan)
```

```
16
17 # plotting the sex covariates using the infants gender , girl and boy
18 kaplan_sex <- survfit(infants_s ~ sex, data = infants)
19 autoplot(kaplan_sex)
20
21 # running the Cox Proportional Hazard Model, using mothers age and infants
        gender as the covariates
22 cox <- coxph(Surv(enter, exit, event) ~ age + sex, data = infants)
23 summary(cox)
```

# Interpretation

```
1
2
3 ############ Interpretation
4
5 # there is a .485 decrease in the expected log of the hazard for male babies
        in comparison to female babies,
6 # while holding the age of the moether constant
7
8 # there is .04 decrease in the expected log of of the hazrd each time the
        mothers age increases by 1 year or 1 unit increase
9 # while holding the sex of infant constant
10
11 # coef    exp(coef) se(coef)                    z Pr(>|z|)
12 # age    -0.04044   0.96037  0.04507 -0.897    0.370
13 # sexboy -0.48518   0.61559  0.44224 -1.097    0.273
14
15
16 ### hazard ratio for male babies is .61 compared to that of female babies or
        62 male babies die for every 100 female babies
17 exp(coef(cox))
18 # age        sexboy
19 # 0.9603673 0.6155879
20
21
22 stargazer(cox)
23
24
25 ######################################### Dependent variable  is : infant_s
26 # testing the out of the coefficients
27 coef(cox_s)
28 # age        sexboy
29 # -0.04043946 -0.48517752
30
31 exp(coef(cox_s))
32 # age     sexboy
33 # 0.9603673 0.6155879
34
35
```
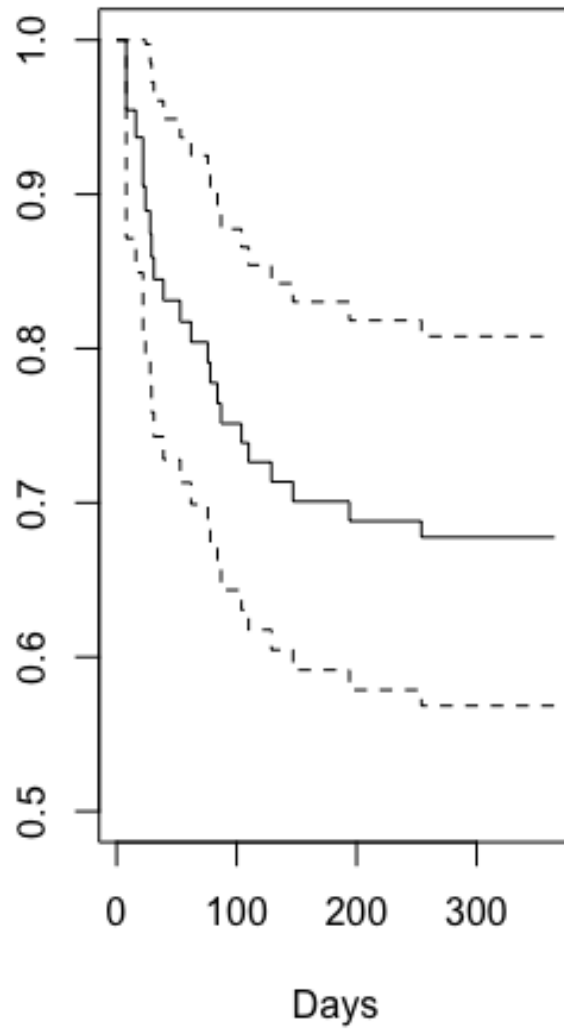
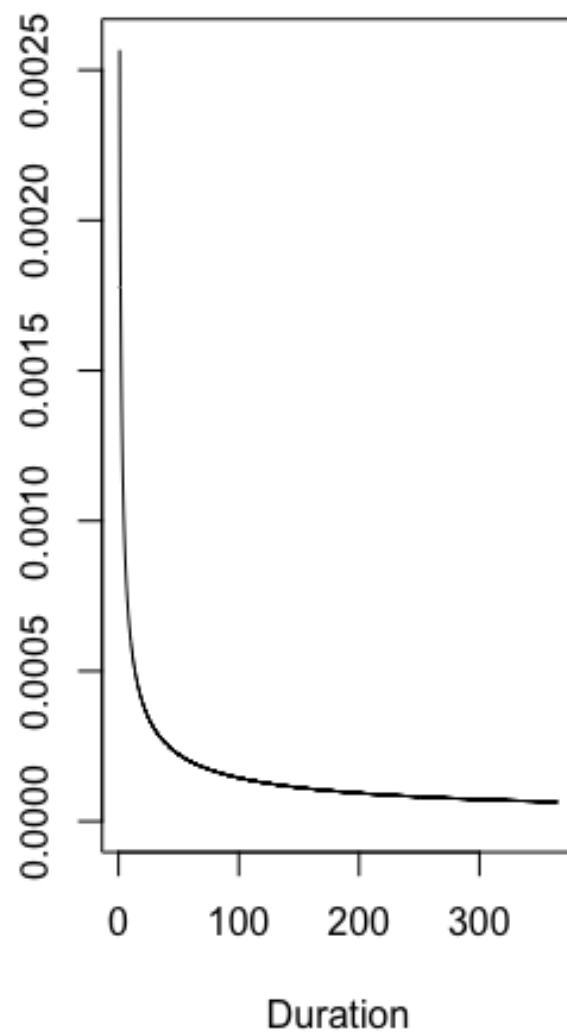Figure 1: Kaplan-Meier Method

**Weibull hazard function**
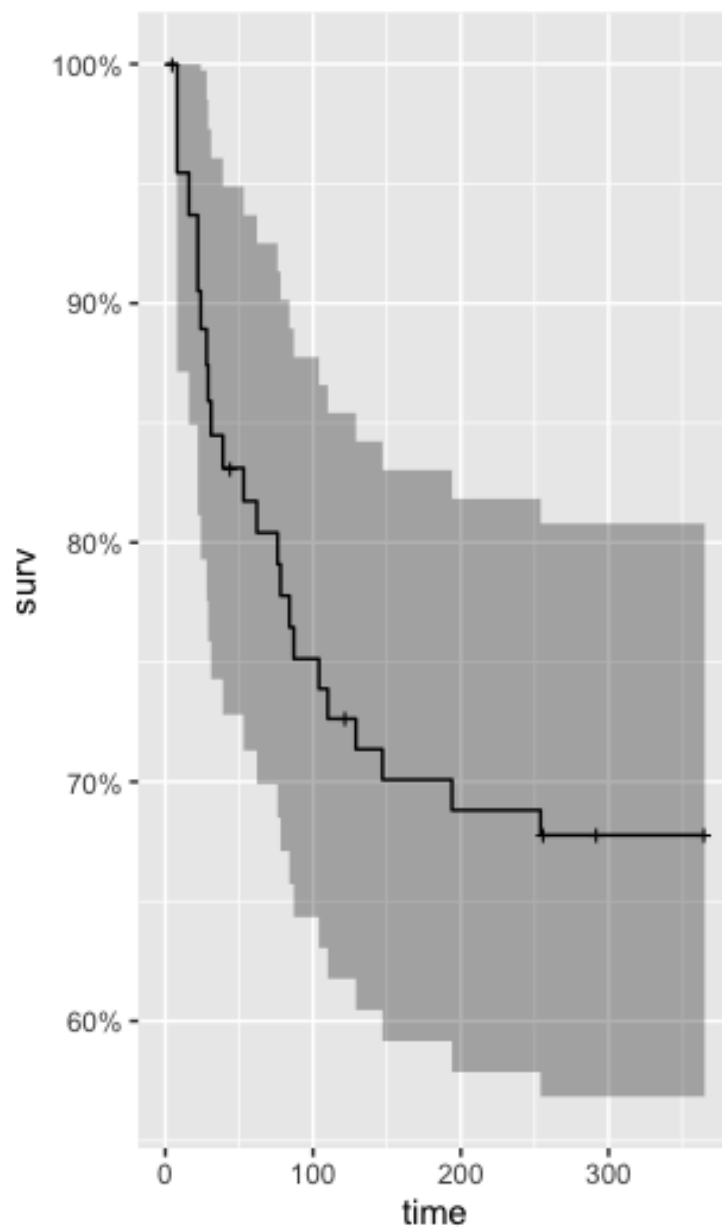


Figure 2: Weibull hazard function
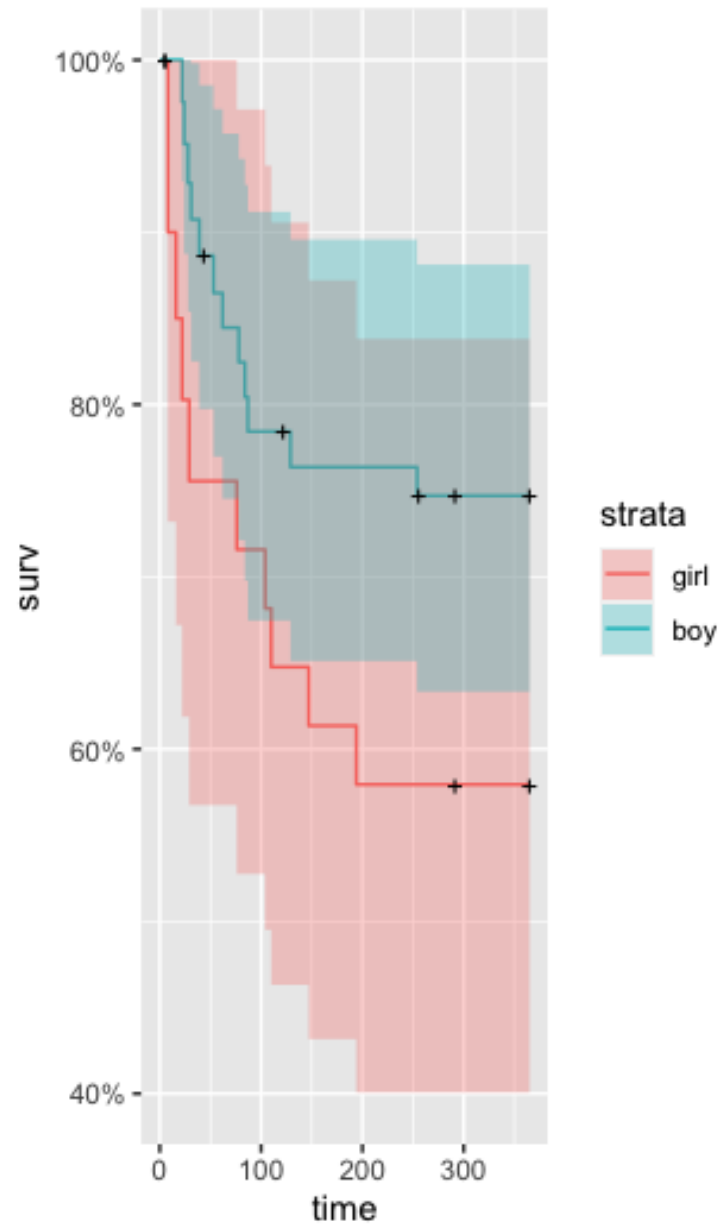
Figure 3: Kaplan-Meir Plot

Figure 4: Kaplan-Meir Plot

```
36 # this command (below) performs a Cox proportional hazards regression analysis
        to estimate the hazard
37 # rate of an event over time, while controlling for the effects of age and sex
        in the "infants" data frame.
38 # This is a common technique used in survival analysis to model the effects of
        predictor variables on
39 # the risk of an event, and to identify factors that may be associated with
        increased or decreased risk.
40
41 cox.w <- phreg(Surv(enter, exit, event) ~ age + sex, data = infants)
42 summary(cox.w)
43 plot(cox.w)
44
45 # Covariate              Mean        Coef      Rel.Risk    S.E.       LR p
46 # age                    27.127     -0.050      0.951      0.045     0.2376
47 # sex                                                                0.4029
48 # girl                   0.317        0          1 (reference)
49 # boy                    0.683      -0.375      0.687      0.444
50
51 # Events                        21
52 # Total time at risk         21616
53 # Max. log. likelihood      -154.85
54 # L3 R test statistic         2.10
```