

Biostat III Examination 2016 Answers

Mark Clements

May 17, 2016

Set-up

```
. global folder 3  
. set linesize 80
```

Commentary

In the following answers, the code and full Stata output are provided together with the answers. The full Stata output was not required in the given answers, but is given here to show how the answers were found.

Some brief comments are warranted on presentation. First, when the question asks for specific results, then those results should be presented separately in text, rather than only presenting the output from the statistical package. Second, the choice of non-proportional fonts makes it difficult to read output from the statistical package. Third, using colours in the graphics makes it difficult to discern which line is which in black-and-white printout. I suggest that using `scheme(s2mono)` would be useful for graphics in Stata.

Part 1

Question 1

We read in the dataset:

```
. import delimited "http://biostat3.net/download/exams/2016/$folder/incidence.c  
> sv", clear  
(6 vars, 360 obs)  
. egen agecat = cut(age), at(40, 50, 60, 70, 80, 90)
```

We then fit a Poisson regression with the number of lung cancer cases at the outcome (first argument), with the person-time of exposure as the `exposure` option. We include attained `age` as a linear, continuous effect in each model.

```
. poisson lc sex age, exposure(pt) nolog irr
```

```
Poisson regression                               Number of obs   =          360  
                                                LR chi2(2)      =        480.18  
                                                Prob > chi2     =          0.0000  
Log likelihood = -839.46342                    Pseudo R2      =          0.2224
```

```
-----+-----  
      lc |           IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]  
-----+-----  
      sex |     2.08734   .1991786     7.71  0.000     1.731289     2.516615  
      age |     1.09231   .0046702    20.65  0.000     1.083195     1.101502  
      _cons |     2.05e-06   5.71e-07   -46.93  0.000     1.18e-06     3.54e-06  
      ln(pt) |           1 (exposure)
```

```
-----
. poisson lc smoking age, exposure(pt) nolog irr
```

```
Poisson regression                                Number of obs   =       360
                                                LR chi2(2)      =      1246.46
                                                Prob > chi2     =       0.0000
Log likelihood = -456.32451                    Pseudo R2       =       0.5773
```

```
-----
      lc |          IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
  smoking |  18.28866   2.285509   23.26  0.000   14.31557   23.36445
    age   |   1.097374   .0047483   21.47  0.000   1.088107   1.10672
  _cons   |  4.71e-07   1.39e-07   -49.15  0.000   2.63e-07   8.41e-07
 ln(pt)   |           1 (exposure)
```

```
-----
. poisson lc asbestos age, exposure(pt) nolog irr
```

```
Poisson regression                                Number of obs   =       360
                                                LR chi2(2)      =       499.97
                                                Prob > chi2     =       0.0000
Log likelihood = -829.57149                    Pseudo R2       =       0.2316
```

```
-----
      lc |          IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
  asbestos |  3.211326   .3645096   10.28  0.000   2.57079   4.011457
    age   |   1.091717   .0046637   20.54  0.000   1.082614   1.100896
  _cons   |  2.77e-06   7.52e-07   -47.17  0.000   1.63e-06   4.71e-06
 ln(pt)   |           1 (exposure)
```

The age-adjusted incidence rate ratio for sex is 2.09 (95% confidence interval (CI): 1.73, 2.52). This association is highly significant ($p < 0.001$).

The age-adjusted incidence rate ratio for smoking is 18.29 (95% confidence interval (CI): 14.32, 23.36). This association is highly significant ($p < 0.001$).

The age-adjusted incidence rate ratio for asbestos is 3.21 (95% confidence interval (CI): 2.57, 4.01). This association is highly significant ($p < 0.001$).

We could have adjusted for attained age in several other ways, including quintiles or splines. To investigate this, we first use quintiles with sex:

```
. xtile ageQ5 = age, nquantiles(5)
. poisson lc sex i.ageQ5, exposure(pt) nolog irr base
```

```
Poisson regression                                Number of obs   =       360
                                                LR chi2(5)      =       463.06
                                                Prob > chi2     =       0.0000
Log likelihood = -848.02325                    Pseudo R2       =       0.2145
```

```
-----
      lc |          IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
    sex   |  2.07351   .1978415    7.64  0.000   1.719846   2.499899
  ageQ5   |
    1     |           1 (base)
    2     |  2.342366   .3938663    5.06  0.000   1.684714   3.256742
    3     |  5.684055   .8926171   11.07  0.000   4.178175   7.732678
    4     | 12.58007   1.972089   16.15  0.000   9.252228  17.10488
```

```

      5 | 15.80467 3.322519 13.13 0.000 10.46749 23.86318
      |
    _cons | .0000916 .0000136 -62.76 0.000 .0000685 .0001225
ln(pt) |          1 (exposure)
-----

```

This shows a very similar point estimate and standard errors to modelling attained age as a linear, continuous effect. We also investigate using restricted cubic splines:

```

. mkspline ageSpline = age, cubic nknots(4)
. poisson lc sex ageSpline*, exposure(pt) nolog irr base

```

```

Poisson regression                               Number of obs   =       360
                                                LR chi2(4)      =       485.97
                                                Prob > chi2     =       0.0000
Log likelihood = -836.5696                    Pseudo R2       =       0.2251
-----

```

```

      lc |          IRR  Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
    sex | 2.081729   .1986431     7.68  0.000    1.726635   2.509849
ageSpline1 | 1.114283   .0236987     5.09  0.000    1.068789   1.161713
ageSpline2 | .9824903   .0612319    -0.28  0.777    .8695183   1.11014
ageSpline3 | .967742    .170512    -0.19  0.852    .6851435   1.366903
    _cons | 7.29e-07   7.60e-07   -13.57  0.000    9.47e-08   5.62e-06
ln(pt) |          1 (exposure)
-----

```

Again, this shows a very similar point estimate and standard errors to modelling attained age as a linear, continuous effect. I accepted answers using any of quintiles, linear/continuous age, splines or similar functional forms.

In summary, lung cancer incidence is associated with age, sex, asbestos exposure and current smoking exposure.

Question 2

We now adjust for age, sex, smoking exposure and asbestos exposure in the same model.

```

. poisson lc age sex smoking asbestos, exposure(pt) nolog irr

```

```

Poisson regression                               Number of obs   =       360
                                                LR chi2(4)      =      1343.01
                                                Prob > chi2     =       0.0000
Log likelihood = -408.05264                    Pseudo R2       =       0.6220
-----

```

```

      lc |          IRR  Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
    age | 1.099109   .00478     21.73  0.000    1.08978   1.108518
    sex | 1.415325   .1366372     3.60  0.000    1.171332   1.710142
 smoking | 17.55656   2.203587    22.83  0.000   13.72783   22.45312
 asbestos | 3.061346   .3502517     9.78  0.000    2.446389   3.830886
    _cons | 3.11e-07   9.48e-08   -49.11  0.000    1.71e-07   5.65e-07
ln(pt) |          1 (exposure)
-----

```

```

. est store ModelA

```

This shows clearly that each of attained age, sex, smoking and asbestos exposure are significantly associated with lung cancer incidence ($p < 0.001$ for all adjusted effects). The adjusted rate ratio (RR)

for age was 1.099 (95% CI: 1.090, 1.109) per year of age, indicating a rapid rise with increasing age. Males have higher rates of disease even after adjustment for other covariates (RR=1.42, 95% CI: 1.17, 1.71). Smoking is strongly associated with lung cancer incidence (RR=17.56, 95% CI: 13.73, 22.45). Finally, asbestos exposure has a rate ratio of 3.06 (95% CI: 2.45, 3.83).

Empirical evidence for confounding can be assessed in several ways. First, we can assess whether exposure to smoking and asbestos are associated:

```
. tab smoking asbestos [aw=pt], row
```

		asbestos		
		0	1	Total
0	253.08739	19.980209		273.0676
	92.68	7.32		100.00
1	80.232149	6.7002514		86.932401
	92.29	7.71		100.00
Total	333.31954	26.6804603		360
	92.59	7.41		100.00

We see that the prevalence of exposure to asbestos is similar or slightly lower among never smokers (7.3%) and current smokers (7.7%). We are not able to undertake a formal statistical test with these weighted and aggregated data.

Second, we can assess whether the estimated associations between lung cancer incidence and each of smoking and asbestos change after an adjustment for other covariates.

Comparing the linear age-adjusted model with the main effects model, we see that the rate ratio for asbestos changed from 3.21 to 3.06 (5% reduction), and the rate ratio for smoking changed from 18.45 to 17.63 (4% reduction). Again, there is limited evidence for confounding between smoking and asbestos.

Question 3

(a)

A regression model formula is

$$\log(\lambda(t|x)) = \beta_0 + \beta_1 \text{age} + \beta_2 I(\text{sex} = 1) + \beta_3 I(\text{smoking} = 1) + \beta_4 I(\text{asbestos} = 1) + \beta_5 I(\text{smoking} = 1 \ \& \ \text{asbestos} = 1)$$

where $\lambda(t|x)$ is the rate at attained age t given covariates x (including sex, smoking and asbestos), with coefficients $\beta_0, \beta_1, \beta_2, \beta_3, \beta_4$ and β_5 , and $I(\text{test})$ is 1 if the test is true and 0 if the test is false.

(b)

We now fit the interaction model:

```
. poisson lc age sex smoking##asbestos, exposure(pt) nolog irr
```

Poisson regression	Number of obs	=	360
	LR chi2(5)	=	1344.29
	Prob > chi2	=	0.0000
Log likelihood = -407.41304	Pseudo R2	=	0.6226

```
-----+-----
```

lc	IRR	Std. Err.	z	P> z	[95% Conf. Interval]	
age	1.099051	.0047794	21.72	0.000	1.089724	1.108459
sex	1.413915	.1364279	3.59	0.000	1.170284	1.708265
1.smoking	18.91174	2.722219	20.42	0.000	14.26286	25.07588
1.asbestos	4.05215	1.07547	5.27	0.000	2.408632	6.817114
smoking# asbestos						
1 1	.7131929	.2091806	-1.15	0.249	.4013731	1.26726
_cons	2.93e-07	9.10e-08	-48.40	0.000	1.59e-07	5.38e-07
ln(pt)	1	(exposure)				

```
-----+-----
```

```
. est store ModelB
. lrtest ModelA ModelB
```

```
Likelihood-ratio test                    LR chi2(1) =      1.28
(Assumption: ModelA nested in ModelB)    Prob > chi2 =      0.2580
```

Comparing Model A with Model B, we see that there is little evidence for a statistical interaction on a multiplicative scale. First, we note that the Wald test for the interaction term has a p-value of 0.25. Second, we see that the likelihood ratio test is also not significant, with $p = 0.26$.

(c)

From Model B, we can calculate the incidence rate for a males aged 62 years who has been exposed to asbestos and is a current smoker using several approaches. We can calculate the rate from the regression estimates, however we need to take account of the covariance terms to calculate the confidence interval, which is best done using tools provided by each statistical package. Using the `lincom` command:

```
. quietly poisson lc age sex smoking##asbestos, exposure(pt) nolog irr
. lincom sex + 1.smoking + 1.asbestos + 1.smoking#1.asbestos + 62*age + _cons,
> irr
```

```
( 1) 62*[lc]age + [lc]sex + [lc]1.smoking + [lc]1.asbestos +
      [lc]1.smoking#1.asbestos + [lc]_cons = 0
```

```
-----+-----
```

lc	IRR	Std. Err.	z	P> z	[95% Conf. Interval]	
(1)	.007901	.0009099	-42.03	0.000	.0063046	.0099017

```
-----+-----
```

This shows that the incidence rate is 7.90 (95% CI: 6.30, 9.90) per 1000 person-years. We could also do this analysis with the `predict` and `margins` command.

Part 2

Question 4

We read in the data using the following:

```
. display "Folder = $folder"
Folder = 3
. import delimited "http://biostat3.net/download/exams/2016/$folder/survival.csv"
> v", clear
(8 vars, 486 obs)
```

(a)

This question is equivalent to completing *Table 1* for a randomised controlled trial to assess whether randomisation led to balanced covariates. We use simple tests to assess whether treatment assignment varies substantially by age at diagnosis, sex, smoking exposure and asbestos exposure.

For age at diagnosis, we can use either a t-test or a non-parametric test:

```
. ttest age, by(tx)
```

Two-sample t test with equal variances

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
0	246	63.33455	.6152337	9.649563	62.12273	64.54638
1	240	62.19994	.6665995	10.32692	60.88678	63.5131
combined	486	62.77425	.4534106	9.995621	61.88336	63.66514
diff		1.134616	.9063606		-.646271	2.915504

```
diff = mean(0) - mean(1)                                t = 1.2518
Ho: diff = 0                                           degrees of freedom = 484
```

```
Ha: diff < 0                Ha: diff != 0                Ha: diff > 0
Pr(T < t) = 0.8944          Pr(|T| > |t|) = 0.2112          Pr(T > t) = 0.1056
. ranksum age, by(tx)
```

Two-sample Wilcoxon rank-sum (Mann-Whitney) test

tx	obs	rank sum	expected
0	246	61571	59901
1	240	56770	58440
combined	486	118341	118341

```
unadjusted variance 2396040.00
adjustment for ties      0.00
-----
adjusted variance 2396040.00
```

```
Ho: age(tx==0) = age(tx==1)
z = 1.079
Prob > |z| = 0.2806
```

We find no evidence that age differs by treatment modality ($p = 0.21$ for the t-test and $p = 0.28$ for the Wilcoxon test). For the other variables:

```
. tab tx sex, chi row
```

```
+-----+
| Key          |
|-----|
| frequency    |
| row percentage |
+-----+

      |          sex
tx |          0          1 | Total
```

	0	1	Total
frequency	88	158	246
row percentage	35.77	64.23	100.00

Pearson chi2(1) = 0.3195 Pr = 0.572
. tab tx smoking, chi row

```
+-----+
| Key      |
|-----|
| frequency|
| row percentage|
+-----+
```

tx	smoking		Total
	0	1	
0	44	202	246
	17.89	82.11	100.00
1	32	208	240
	13.33	86.67	100.00
Total	76	410	486
	15.64	84.36	100.00

Pearson chi2(1) = 1.9088 Pr = 0.167
. tab tx asbestos, chi row

```
+-----+
| Key      |
|-----|
| frequency|
| row percentage|
+-----+
```

tx	asbestos		Total
	0	1	
0	194	52	246
	78.86	21.14	100.00
1	195	45	240
	81.25	18.75	100.00
Total	389	97	486
	80.04	19.96	100.00

Pearson chi2(1) = 0.4337 Pr = 0.510

We find little evidence that randomisation varied by sex ($p = 0.57$), by smoking ($p = 0.17$) or by asbestos exposure ($p = 0.51$). We could check for potential confounding by sex in the survival analysis.

(b)

We `stset` the data using time since diagnosis as the primary time scale and then plot the Kaplan-Meier curves

```
. stset tsurv, failure(event) id(id)

           id: id
failure event: event != 0 & event < .
obs. time interval: (tsurv[_n-1], tsurv]
exit on or before: failure

-----
486 total observations
  0 exclusions
-----

486 observations remaining, representing
486 subjects
424 failures in single-failure-per-subject data
583.8291 total analysis time at risk and under observation
                                at risk from t =          0
                                earliest observed entry t =      0
                                last observed exit t =          5
. sts graph, by(tx) name(km1, replace) scheme(s2mono)

           failure _d: event
           analysis time _t: tsurv
           id: id
. graph export exam_2016_km1.eps, name(km1) replace
(file exam_2016_km1.eps written in EPS format)
. * the following line is only needed on Linux
. !! convert -density 300 exam_2016_km1.eps exam_2016_km1_$folder.png
. sts test tx

           failure _d: event
           analysis time _t: tsurv
           id: id
```

Log-rank test for equality of survivor functions

tx	Events observed	Events expected
0	210	239.11
1	214	184.89
Total	424	424.00

```
chi2(1) =      8.16
Pr>chi2 =     0.0043
```

```
. sts list, by(tx) at(1 2 3 4 5)
```

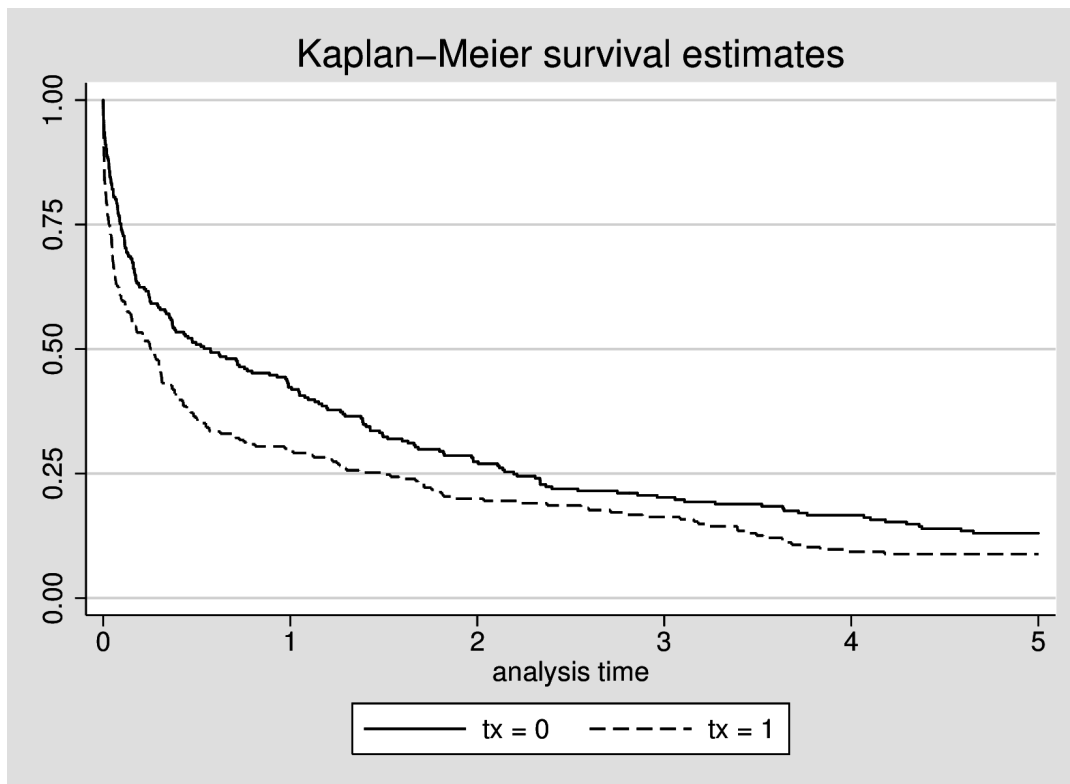
```
           failure _d: event
           analysis time _t: tsurv
           id: id
```

Beg. Survivor Std.

Time	Total	Fail	Function	Error	[95% Conf. Int.]	

tx=0						
1	104	141	0.4231	0.0316	0.3607	0.4842
2	67	36	0.2739	0.0286	0.2194	0.3311
3	47	17	0.2021	0.0259	0.1539	0.2549
4	38	8	0.1663	0.0242	0.1221	0.2165
5	29	8	0.1303	0.0221	0.0909	0.1770
tx=1						
1	70	167	0.3001	0.0298	0.2431	0.3591
2	45	23	0.1997	0.0262	0.1511	0.2531
3	36	8	0.1628	0.0244	0.1184	0.2134
4	21	15	0.0930	0.0195	0.0594	0.1356
5	18	1	0.0884	0.0190	0.0557	0.1302

Note: survivor function is calculated over full data and evaluated at indicated times; it is not calculated from aggregates shown at left.



The Kaplan-Meier curves show that survival is poor for lung cancer patients, with fewer than 25% of patients surviving to 5 years. We also see that treatment with chemotherapy+radiotherapy leads to more deaths soon after diagnosis. It is unclear whether the rates or hazards are different after one year.

Although not specifically asked for, we also (i) used the log-rank test to compare the curves, finding strong evidence for a difference ($p = 0.004$) and (ii) estimated survival to five years, where 13% (95% CI: 9, 18) survived for those on conventional treatment and 9% (95% CI: 6, 13) survived for those on chemotherapy+radiotherapy.

Question 5

Based on Question 4 (a), we first investigated whether age and sex were associated with survival and hence would be potential confounders:

```
. stcox tx sex age, nolog
      failure _d: event
```

```
analysis time _t: tsurv
id: id
```

Cox regression -- no ties

```
No. of subjects = 486          Number of obs = 486
No. of failures = 424
Time at risk = 583.8291199
Log likelihood = -2307.7533    LR chi2(3) = 11.12
                               Prob > chi2 = 0.0111
```

_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
tx	1.310947	.1278614	2.78	0.006	1.08284	1.587106
sex	1.129684	.117439	1.17	0.241	.9214424	1.384987
age	.9947913	.0048869	-1.06	0.288	.9852591	1.004416

. stcox tx sex, nolog

```
failure _d: event
analysis time _t: tsurv
id: id
```

Cox regression -- no ties

```
No. of subjects = 486          Number of obs = 486
No. of failures = 424
Time at risk = 583.8291199
Log likelihood = -2308.3165    LR chi2(2) = 9.99
                               Prob > chi2 = 0.0068
```

_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
tx	1.319648	.1284444	2.85	0.004	1.090458	1.59701
sex	1.150779	.1179678	1.37	0.171	.9413129	1.406856

. stcox tx age, nolog

```
failure _d: event
analysis time _t: tsurv
id: id
```

Cox regression -- no ties

```
No. of subjects = 486          Number of obs = 486
No. of failures = 424
Time at risk = 583.8291199
Log likelihood = -2308.4493    LR chi2(2) = 9.72
                               Prob > chi2 = 0.0077
```

_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
tx	1.309927	.1277078	2.77	0.006	1.082085	1.585742
age	.9938174	.0048068	-1.28	0.200	.9844407	1.003283

```

. stcox tx, nolog

      failure _d:  event
analysis time _t:  tsurv
              id:  id

Cox regression -- no ties

No. of subjects =          486          Number of obs =          486
No. of failures =          424
Time at risk    = 583.8291199

Log likelihood   = -2309.2681          LR chi2(1)      =          8.09
                                          Prob > chi2    =          0.0045

```

```

-----+-----
      _t | Haz. Ratio   Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
      tx |   1.319277   .1283979    2.85   0.004    1.090167    1.596536
-----+-----

```

Adjusting for treatment modality, there is no evidence that either sex or age are associated with survival, with Wald test p-values of 0.17 and 0.20 for sex and age, respectively. Furthermore, fitting a Cox regression models with and without age and sex suggest that the effect of treatment modality is insensitive to inclusion of age and sex in the model. The hazard ratio for chemotherapy+radiotherapy compared with conventional therapy is 1.32 (95% CI: 1.09, 1.60), suggesting that the average hazard ratio for chemotherapy+radiotherapy is high over the five-year period.

For the time scale, we have initially used time since cancer diagnosis. There is a strong association between time since diagnosis and survival, suggesting that this is the best choice of primary time scale. Moreover, there is a suggestion of non-proportional hazards, with a higher rate ratio in the first year than for the later years. We could investigate using attained age as the primary time scale, but then we would need to finely model for the time since diagnosis, which would require modelling two time scales. For simplicity, we propose using time since diagnosis as the primary time scale.

Question 6

(i)

For an analysis of scaled Schoenfeld residuals, we use:

```

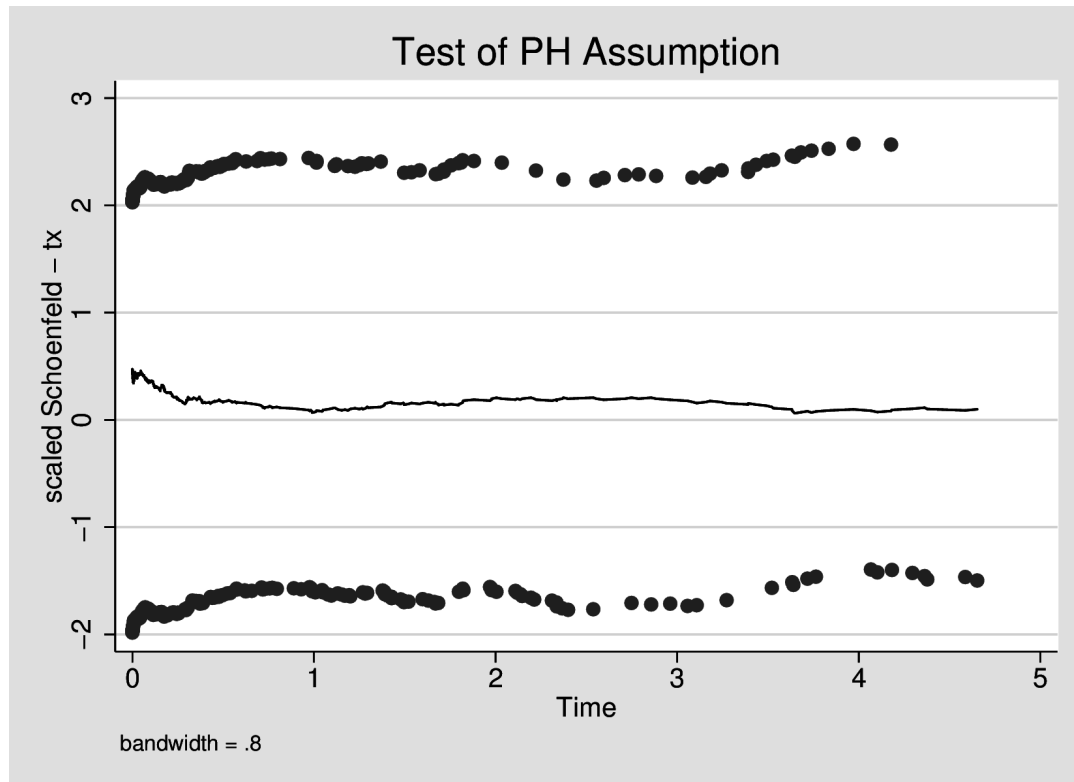
. estat phtest, detail

Test of proportional-hazards assumption

Time:  Time
-----+-----
      |      rho      chi2      df      Prob>chi2
-----+-----
      tx      |   -0.06664    1.85      1      0.1737
-----+-----
global test |          1.85      1      0.1737
-----+-----

. estat phtest, plot(tx) name(phtest, replace) scheme(s2mono)
. graph export exam_2016_phtest.eps, name(phtest) replace
(file exam_2016_phtest.eps written in EPS format)
. * the following line is only needed on Linux
. !! convert -density 300 exam_2016_phtest.eps exam_2016_phtest_$folder.png

```



This shows that there is little evidence ($p = 0.17$) that the hazard ratio decreases with increasing time since diagnosis: the scaled residuals and linear time have a correlation of -0.07 . From the plot of the scaled residuals and time, we see the running mean smoother dips early in the follow-up period and then is flat or very slightly declining. Given the number of events that are early in the period, we could also test using a log-transformation for time since diagnosis:

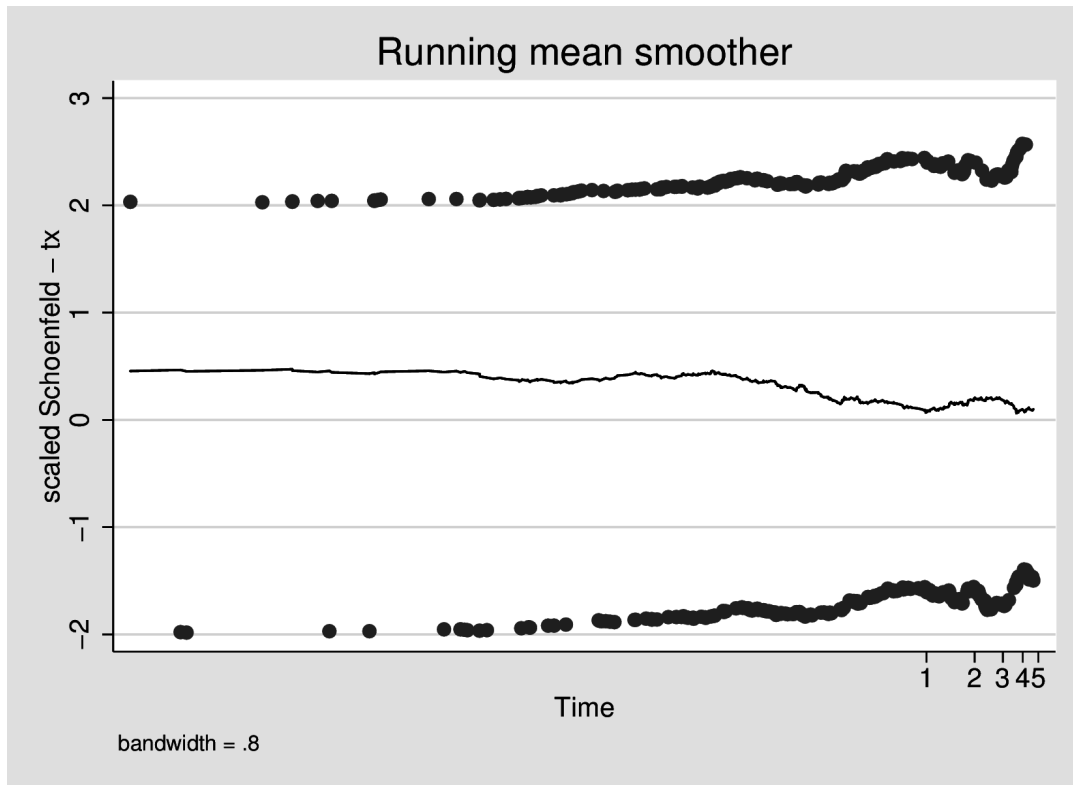
```
. estat phtest, detail log
```

```
Test of proportional-hazards assumption
```

```
Time: Log(t)
```

	rho	chi2	df	Prob>chi2
tx	-0.10910	4.96	1	0.0259
global test		4.96	1	0.0259

```
. estat phtest, log plot(tx) name(phtestlog, replace) scheme(s2mono)
. graph export exam_2016_phtestlog.eps, name(phtestlog) replace
(file exam_2016_phtestlog.eps written in EPS format)
. * the following line is only needed on Linux
. !! convert -density 300 exam_2016_phtestlog.eps exam_2016_phtestlog_$folder.p
> ng
```



In contrast to the linear time scale, there is evidence for non-proportionality on a log(time) scale ($p = 0.03$).

(ii)

We can test for piecewise-constant hazard ratios by splitting by time and fitting for an interaction. In the following, the "c" prefix indicates a continuous variable, while the "i" prefix indicates a factor variable.

```
. quietly import delimited "http://biostat3.net/download/exams/2016/$folder/sur
> vival.csv", clear
. quietly stset tsurv, fail(event) id(id)
. stsplit timeband, at(0, 1, max)
(172 observations (episodes) created)
. stcox sex i.tx##i.timeband, nolog
```

```
      failure _d:  event
analysis time _t:  tsurv
              id:  id
```

Cox regression -- no ties

```
No. of subjects =          486                Number of obs =          658
No. of failures =          424
Time at risk    = 583.8291199
Log likelihood  = -2306.7568                LR chi2(3) =          13.11
                                                Prob > chi2 =          0.0044
```

	_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
	sex	1.149816	.1178737	1.36	0.173	.9405183 1.40569
	1.tx	1.466359	.1679591	3.34	0.001	1.1715 1.835432
	1.timeband	20.08552

```

      |
tx#timeband |
      1 1 | .6780242 .1500214 -1.76 0.079 .4394466 1.046127
-----

```

```
. stcox tx sex c.tx#c.timeband, nolog
```

```

      failure _d: event
analysis time _t: tsurv
              id: id

```

```
Cox regression -- no ties
```

```

No. of subjects =          486          Number of obs =          658
No. of failures =          424
Time at risk    = 583.8291199
Log likelihood   = -2306.7568          LR chi2(3)    =          13.11
                                          Prob > chi2   =          0.0044

```

```

-----
      _t | Haz. Ratio   Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
      tx |  1.466359   .1679591     3.34  0.001     1.1715   1.835432
      sex |  1.149816   .1178737     1.36  0.173     .9405183  1.40569
      |
      c.tx# |
c.timeband | .6780242   .1500214    -1.76  0.079     .4394466  1.046127
-----

```

```
. stcox c.tx#i.timeband, nolog
```

```

      failure _d: event
analysis time _t: tsurv
              id: id

```

```
Cox regression -- no ties
```

```

No. of subjects =          486          Number of obs =          658
No. of failures =          424
Time at risk    = 583.8291199
Log likelihood   = -2307.697          LR chi2(2)    =          11.23
                                          Prob > chi2   =          0.0036

```

```

-----
      _t | Haz. Ratio   Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
timeband#c.tx |
      0 |  1.466498   .1679618     3.34  0.001     1.171632  1.835574
      1 |  .9929335   .1879615    -0.04  0.970     .6851542  1.438971
-----

```

This model provides some evidence that the hazard ratio is time-dependent for a piecewise-constant hazard ratio ($p = 0.08$). The hazard ratio in the first year is raised at 1.47 (95% CI: 1.17, 1.84), while the hazard ratio after the first year is close to 1 (HR=0.99; 95% CI: 0.69, 1.44).

(iii)

We can re-fit the model in (ii) using Stata `stcox`'s `tv` and `te` options:

```
. stcox tx, nolog tv(c.tx) te(_t>=1)
```

```

failure _d: event
analysis time _t: tsurv
id: id

```

Cox regression -- no ties

```

No. of subjects =          486          Number of obs =          658
No. of failures =          424
Time at risk    = 583.8291199
Log likelihood  = -2307.697          LR chi2(2) =          11.23
                                          Prob > chi2 =          0.0036

```

	_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
-----+-----						
main						
	tx	1.466498	.1679618	3.34	0.001	1.171632 1.835574
-----+-----						
tvc						
	tx	.6770779	.149804	-1.76	0.078	.4388435 1.044642
-----+-----						

Note: variables in tvc equation interacted with _t>=1

Again, we find some evidence for a time-dependent hazard ratio ($p = 0.08$). We can model for a time-dependent hazard ratio that depends on time:

```
. stcox tx, nolog tvc(tx) texp(_t)
```

```

failure _d: event
analysis time _t: tsurv
id: id

```

Cox regression -- no ties

```

No. of subjects =          486          Number of obs =          658
No. of failures =          424
Time at risk    = 583.8291199
Log likelihood  = -2308.3273          LR chi2(2) =           9.97
                                          Prob > chi2 =          0.0068

```

	_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
-----+-----						
main						
	tx	1.449659	.1731819	3.11	0.002	1.147039 1.832119
-----+-----						
tvc						
	tx	.8829756	.0808914	-1.36	0.174	.7378503 1.056645
-----+-----						

Note: variables in tvc equation interacted with _t

The interpretation of this model is as follows: the hazard ratio at time 0 is 1.45 (95% CI: 1.15, 1.83). For every year, there is little evidence for a linear decrease in the hazard ratio (RR=0.99, 95% CI: 0.98, 1.00).

(iv)

Using `stpm2` with time-dependent hazard ratios, we use a low-dimensional natural spline for the time-dependent effect. We use a Wald test to check for time-dependence and plot the time-dependent hazard ratio:

```
. stpm2 tx, df(4) scale(hazard) nolog eform tvc(tx) dftvc(2)
note: delayed entry models are being fitted
```

```
Log likelihood = -1084.5605                Number of obs   =          658
```

```
-----+-----
```

	exp(b)	Std. Err.	z	P> z	[95% Conf. Interval]
xb					
tx	1.478171	.1648009	3.51	0.000	1.188022 1.839182
_rcs1	3.289988	.2999739	13.06	0.000	2.75159 3.933734
_rcs2	1.06185	.0695569	0.92	0.360	.9339096 1.207317
_rcs3	.9966186	.0238257	-0.14	0.887	.9509982 1.044427
_rcs4	.9770622	.0151503	-1.50	0.135	.9478149 1.007212
_rcs_tx1	.8450199	.09426	-1.51	0.131	.6790744 1.051518
_rcs_tx2	1.025035	.0771313	0.33	0.742	.8844803 1.187927
_cons	.5187743	.0435441	-7.82	0.000	.4400799 .6115406

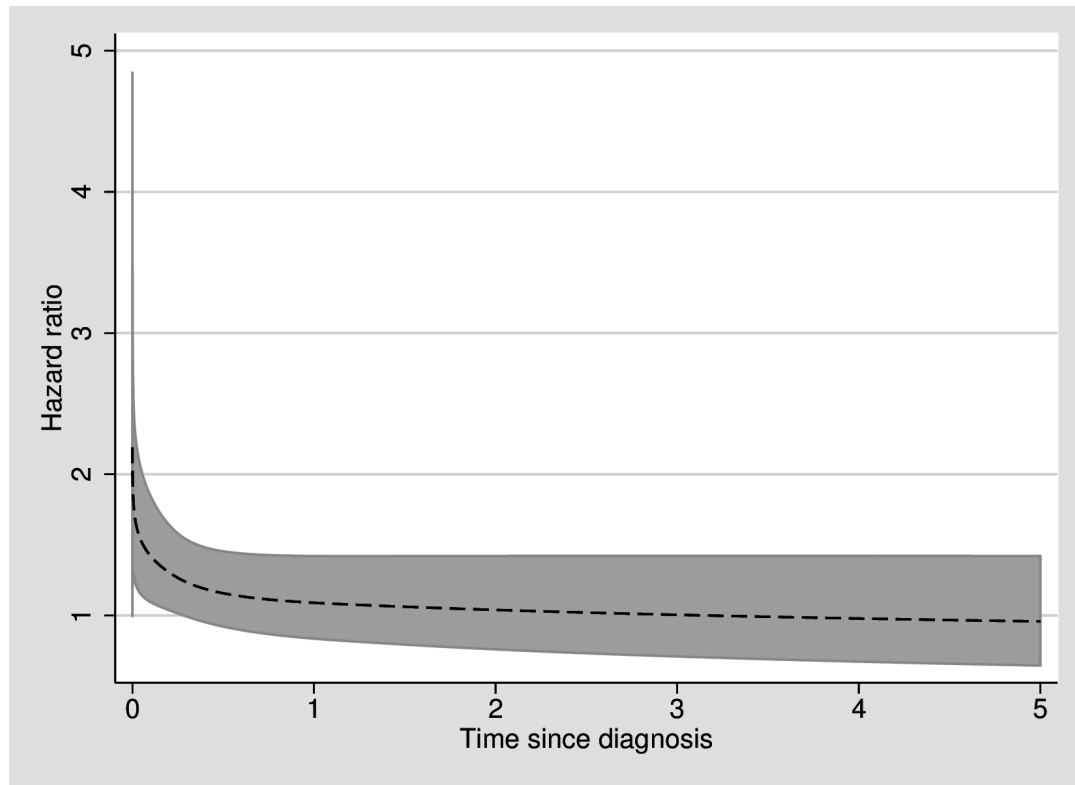
```
-----+-----
```

```
. test _rcs_tx1 _rcs_tx2
```

```
( 1) [xb]_rcs_tx1 = 0
( 2) [xb]_rcs_tx2 = 0
```

```
chi2( 2) = 4.95
Prob > chi2 = 0.0843
```

```
. predict hr, hrnumerator(tx 1) ci
. twoway (rarea hr_lci hr_uci _t if hr_uci<5, sort color(gs12)) (line hr _t if
> hr_uci<5, sort), legend(off) xtitle("Time since diagnosis") ytitle("Hazard ra
> tio") name(hr, replace) scheme(s2mono)
. graph export exam_2016_hr.eps, name(hr) replace
(file exam_2016_hr.eps written in EPS format)
. * the following line is only needed on Linux
. !! convert -density 300 exam_2016_hr.eps exam_2016_hr_$folder.png
```

We see that there is some evidence for time-dependent hazards ($p = 0.08$ from the Wald test). We also see from the plot that the hazard ratio looks comparatively stable across the follow-up period.

Question 7

(a)

Advantages of using Poisson regression for Questions 5–6 include: (i) Poisson regression readily models for multiple time scales, where we could split on attained age and time since diagnosis and then model for main effects and interactions between those time scales and interactions between a time scale and other covariates; (ii) it is simpler to predict rates from Poisson regression, as the analysis is done on that scale.

Disadvantages of using Poisson regression include: (i) the need to split on the time scales, which may increase the size of the computational problem; (ii) the need to specify a functional form for the primary time scale using parametric functions, rather than using Cox regression's non-parametric formulation; (iii) crude time splitting will assume that rates are piece-wise constant, which may not be appropriate; (iv) risk calculations for Poisson regression require that the risk period involves constant rates or numerical integration.

(b)

Assuming that the follow-up time has been split for within one year of diagnosis and from one year of diagnosis, we can model the rate using:

$$\log(\lambda(t|\text{tx})) = \beta_0 + \beta_1 I(t < 1) + \beta_2 I(t \geq 1) + \beta_3 I(\text{tx} = 1) + \beta_4 I(\text{tx} = 1 \ \& \ t \geq 1)$$

A better formulation would be to include more time-splits for time since diagnosis. If we let time cuts be represented by t_j where $t_0 = 0$, then

$$\log(\lambda(t|\text{tx})) = \beta_0 + \sum_j \beta_j I(t_{j-1} < t \leq t_j) + \beta_{\text{tx}} I(\text{tx} = 1) + \beta_{\text{tx},t} I(\text{tx} = 1 \ \& \ t \geq 1)$$

We could also model using splines. Any similar formulation was accepted, including different formulations for the time-dependent hazard ratios.