

Biostat III Examination 2016 Answers

Mark Clements

May 17, 2016

Set-up

```
. global folder 4  
. 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)      =          509.71  
                                                Prob > chi2     =           0.0000  
Log likelihood = -839.02839                    Pseudo R2       =           0.2330
```

```
-----+-----  
      lc |           IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]  
-----+-----  
      sex |    2.336099   .2192283     9.04  0.000     1.94362   2.807833  
      age |    1.089798   .0044831    20.90  0.000     1.081047   1.098621  
      _cons |    2.33e-06   6.27e-07   -48.18  0.000     1.37e-06   3.95e-06  
      ln(pt) |           1   (exposure)
```

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

```
Poisson regression                                Number of obs   =       360
                                                    LR chi2(2)      =      1229.43
                                                    Prob > chi2     =       0.0000
Log likelihood = -479.16733                       Pseudo R2      =       0.5620
```

```
-----
      lc |          IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
  smoking |   15.00551   1.691231    24.03   0.000    12.03134    18.7149
    age   |    1.095422   .004574    21.83   0.000    1.086493    1.104424
  _cons   |   6.63e-07   1.87e-07   -50.34   0.000    3.81e-07    1.15e-06
 ln(pt)   |             1 (exposure)
```

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

```
Poisson regression                                Number of obs   =       360
                                                    LR chi2(2)      =       503.83
                                                    Prob > chi2     =       0.0000
Log likelihood = -841.97002                       Pseudo R2      =       0.2303
```

```
-----
      lc |          IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
  asbestos |   3.088896   .3411285    10.21   0.000    2.487707    3.835372
    age    |    1.088151   .0044504    20.66   0.000    1.079463    1.096909
  _cons    |   3.62e-06   9.37e-07   -48.43   0.000    2.18e-06    6.01e-06
 ln(pt)    |             1 (exposure)
```

The age-adjusted incidence rate ratio for sex is 2.16 (95% confidence interval (CI): 1.80, 2.60). This association is highly significant ($p < 0.001$).

The age-adjusted incidence rate ratio for smoking is 18.45 (95% confidence interval (CI): 14.56, 23.37). This association is highly significant ($p < 0.001$).

The age-adjusted incidence rate ratio for asbestos is 3.68 (95% confidence interval (CI): 2.99, 4.53). 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)      =       493.26
                                                    Prob > chi2     =       0.0000
Log likelihood = -847.25205                       Pseudo R2      =       0.2255
```

```
-----
      lc |          IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
    sex   |   2.315037   .2172295     8.95   0.000    1.926133    2.782465
  ageQ5  |
    1     |             1 (base)
    2     |   2.804366   .456974     6.33   0.000    2.037654    3.859568
    3     |   6.086247   .947346    11.60   0.000    4.485977    8.257377
    4     |  12.54986    1.967415    16.14   0.000    9.229916   17.06396
```

```

      5 | 17.41423 3.532478 14.09 0.000 11.70141 25.91616
      |
    _cons | .0000849 .0000126 -63.03 0.000 .0000634 .0001136
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)      =          517.12
                                                Prob > chi2     =           0.0000
Log likelihood = -835.32277                    Pseudo R2       =           0.2364
-----

```

```

      lc |          IRR  Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
      sex | 2.325112   .2182051     8.99  0.000    1.934465   2.794645
ageSpline1 | 1.117549   .0228575     5.43  0.000    1.073635   1.163259
ageSpline2 | .9694936   .0582494    -0.52  0.606    .8617927   1.090654
ageSpline3 | .997276    .1698155    -0.02  0.987    .7142879   1.392379
    _cons | 6.46e-07   6.47e-07   -14.22  0.000    9.06e-08   4.60e-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)      =         1343.33
                                                Prob > chi2     =           0.0000
Log likelihood = -422.22001                    Pseudo R2       =           0.6140
-----

```

```

      lc |          IRR  Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
      age | 1.097571   .0046186    22.12  0.000    1.088556   1.106661
      sex | 1.58053    .1507423     4.80  0.000    1.311052   1.905398
    smoking | 14.37208   1.629092    23.51  0.000   11.50893   17.94752
    asbestos | 3.014193   .3365934     9.88  0.000    2.421686   3.751668
    _cons | 4.01e-07   1.18e-07   -50.15  0.000    2.25e-07   7.13e-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 1.104 (95% CI: 1.095, 1.113) per year of age, indicating a rapid rise with increasing age. Males have higher rates of disease even after adjustment for other covariates (RR=1.45, 95% CI: 1.20, 1.74). Smoking is strongly associated with lung cancer incidence (RR=17.63, 95% CI: 13.90, 22.35). Finally, asbestos exposure has a rate ratio of 3.27 (95% CI: 2.64, 4.05).

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
```

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

      |          asbestos
      |          0          1 |          Total
-----+-----+-----
      0 | 252.81839  20.275815 | 273.09421
      |          92.58      7.42 |          100.00
-----+-----+-----
      1 | 80.590127  6.3156648 | 86.905792
      |          92.73      7.27 |          100.00
-----+-----+-----
      Total | 333.40852  26.59148 |          360
      |          92.61      7.39 |          100.00
```

We see that the prevalence of exposure to asbestos is similar or slightly lower among never smokers (7.4%) and current smokers (8.0%). We are not able to undertake a formal statistical test with these weighted 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 3.68 to 3.42 (7% 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)         =        1343.86
                             Prob > chi2        =          0.0000
Log likelihood = -421.95474   Pseudo R2         =          0.6143
```

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

lc	IRR	Std. Err.	z	P> z	[95% Conf. Interval]	
age	1.097507	.0046182	22.11	0.000	1.088493	1.106596
sex	1.579898	.1506106	4.80	0.000	1.310644	1.904467
1.smoking	15.00148	1.932057	21.03	0.000	11.65485	19.30906
1.asbestos	3.528549	.8446109	5.27	0.000	2.207234	5.640842
smoking# asbestos						
1 1	.8201726	.2209686	-0.74	0.462	.4837009	1.390701
_cons	3.88e-07	1.15e-07	-49.64	0.000	2.17e-07	6.96e-07
ln(pt)	1	(exposure)				

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

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

```
Likelihood-ratio test                    LR chi2(1) =    0.53
(Assumption: ModelA nested in ModelB)    Prob > chi2 =    0.4664
```

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 0.18. Second, we see that the likelihood ratio test is also not significant, with $p = 0.19$.

(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)	.0085281	.0009697	-41.90	0.000	.0068244	.0106572

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

This shows that the incidence rate is 9.19 (95% CI: 7.46, 10.32) per 1000 person-years.

We can do the same analysis using the `margins` command:

```
. margins smoking##asbestos, predict(ir) at(age=62 sex=1)
```

```
Predictive margins                    Number of obs   =    360
Model VCE      : OIM
```

```
Expression   : Predicted incidence rate, predict(ir)
at           : age           =    62
              sex           =    1
```

		Delta-method				[95% Conf. Interval]	
		Margin	Std. Err.	z	P> z		
smoking							
	0	.0004448	.0000744	5.98	0.000	.0002989	.0005906
	1	.0057375	.0005	11.48	0.000	.0047576	.0067174
asbestos							
	0	.0015716	.0000951	16.52	0.000	.0013852	.0017581
	1	.0046106	.0004915	9.38	0.000	.0036473	.0055739
smoking#							
asbestos							
	0 0	.0001964	.0000244	8.05	0.000	.0001486	.0002443
	0 1	.0006931	.0001457	4.76	0.000	.0004075	.0009787
	1 0	.0029468	.0001844	15.98	0.000	.0025854	.0033082
	1 1	.0085281	.0009697	8.79	0.000	.0066275	.0104288

Finally, we could also do this analysis with the `predict` command.

Part 2

Question 4

We read in the data using the following:

```
. display "Folder = $folder"
Folder = 4
. import delimited "http://biostat3.net/download/exams/2016/$folder/survival.csv", clear
(8 vars, 520 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	254	62.42431	.6415783	10.22508	61.16079	63.68782
1	266	62.7812	.6013364	9.807499	61.59719	63.9652
combined	520	62.60687	.4387732	10.00557	61.74488	63.46886
diff		-.3568897	.8784871		-2.082725	1.368946

```
diff = mean(0) - mean(1)                                t = -0.4063
Ho: diff = 0                                           degrees of freedom = 518
```

```
Ha: diff < 0                                           Ha: diff != 0                                           Ha: diff > 0
Pr(T < t) = 0.3424                                     Pr(|T| > |t|) = 0.6847                                   Pr(T > t) = 0.6576
. ranksum age, by(tx)
```

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

tx	obs	rank sum	expected
0	254	65485	66167
1	266	69975	69293
combined	520	135460	135460

unadjusted variance 2933403.67
 adjustment for ties 0.00

 adjusted variance 2933403.67

Ho: age(tx==0) = age(tx==1)
 z = -0.398
 Prob > |z| = 0.6905

We find no evidence that age differs by treatment modality ($p = 0.46$ for the t-test and $p = 0.61$ for the Wilcoxon test). For the other variables:

. tab tx sex, chi row

tx	sex		Total
	0	1	
0	92	162	254
	36.22	63.78	100.00
1	76	190	266
	28.57	71.43	100.00
Total	168	352	520
	32.31	67.69	100.00

Pearson chi2(1) = 3.4760 Pr = 0.062

. tab tx smoking, chi row

tx	smoking		Total
	0	1	
0	47	207	254
	18.50	81.50	100.00

1	50	216	266
	18.80	81.20	100.00
-----+			
Total	97	423	520
	18.65	81.35	100.00

Pearson chi2(1) = 0.0074 Pr = 0.932

. tab tx asbestos, chi row

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

tx	asbestos		Total
	0	1	
0	211	43	254
	83.07	16.93	100.00
1	207	59	266
	77.82	22.18	100.00
Total	418	102	520
	80.38	19.62	100.00

Pearson chi2(1) = 2.2724 Pr = 0.132

We find little evidence that randomisation varied by sex ($p = 0.09$), by smoking ($p = 0.21$) or by asbestos exposure ($p = 0.86$). 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
```

```
-----
520 total observations
  0 exclusions
-----
```

```
520 observations remaining, representing
520 subjects
461 failures in single-failure-per-subject data
528.6148 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	221	255.26
1	240	205.74
Total	461	461.00

chi2(1) = 10.36

Pr>chi2 = 0.0013

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

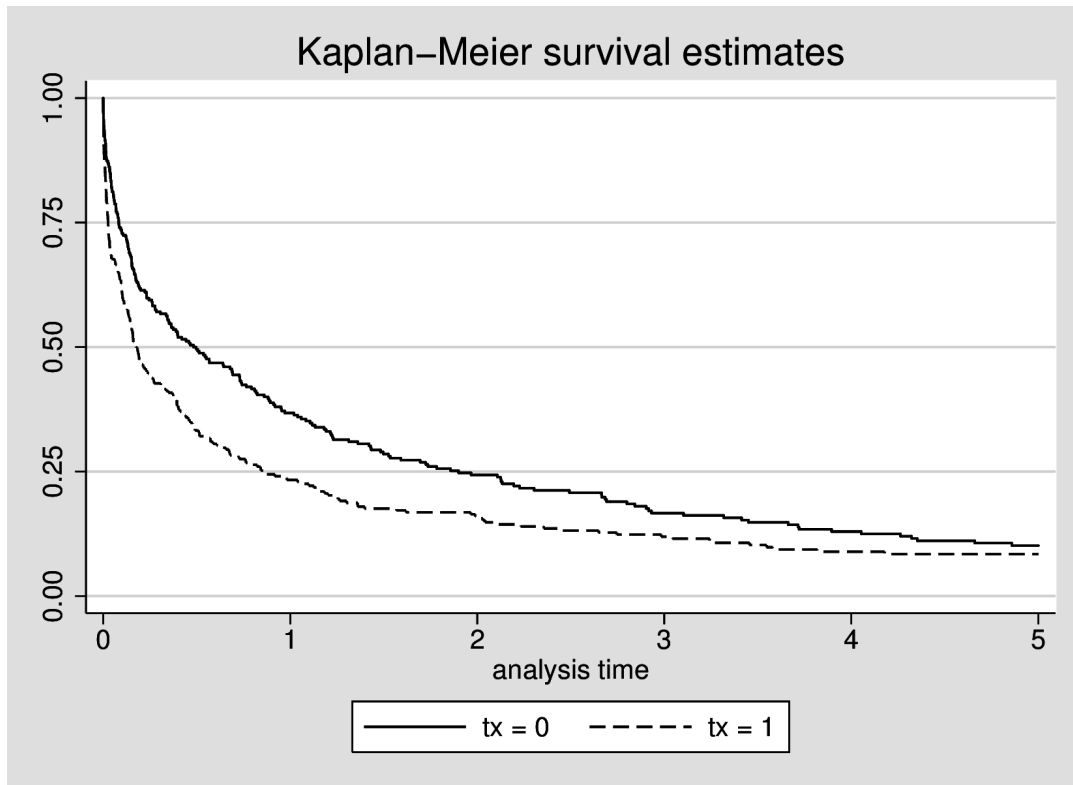
```

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

```

	Time	Beg. Total	Fail	Survivor Function	Std. Error	[95% Conf. Int.]	
tx=0							
	1	91	160	0.3677	0.0304	0.3085	0.4270
	2	57	30	0.2432	0.0273	0.1917	0.2982
	3	37	17	0.1665	0.0242	0.1222	0.2168
	4	29	8	0.1295	0.0221	0.0900	0.1763
	5	21	6	0.1015	0.0201	0.0666	0.1450
tx=1							
	1	62	203	0.2331	0.0260	0.1840	0.2856
	2	41	19	0.1600	0.0227	0.1185	0.2070
	3	30	10	0.1193	0.0202	0.0832	0.1622
	4	21	7	0.0891	0.0181	0.0579	0.1285
	5	16	1	0.0842	0.0177	0.0537	0.1231

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 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 ($p = 0.0001$) and (ii) estimated survival to five years, where 9% (95% CI: 6, 13) survived for those on conventional treatment and 3% (95% CI: 1, 6) 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 =          520          Number of obs   =          520
No. of failures =          461
Time at risk    = 528.6147698
Log likelihood  = -2530.0625          LR chi2(3)       =          10.34
                                          Prob > chi2     =          0.0159
```

	_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
	tx	1.346988	.1263471	3.18	0.001	1.120782 1.618848
	sex	1.023571	.1020758	0.23	0.815	.8418439 1.244526
	age	1.000234	.0046747	0.05	0.960	.9911134 1.009438

Adjusting for treatment modality, there is no evidence that either sex or age are associated with survival, with Wald test p-values of 0.75 and 0.25 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 1.77 (95% CI: 1.47, 2.13), 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	-0.12839	7.45	1	0.0063
global test		7.45	1	0.0063

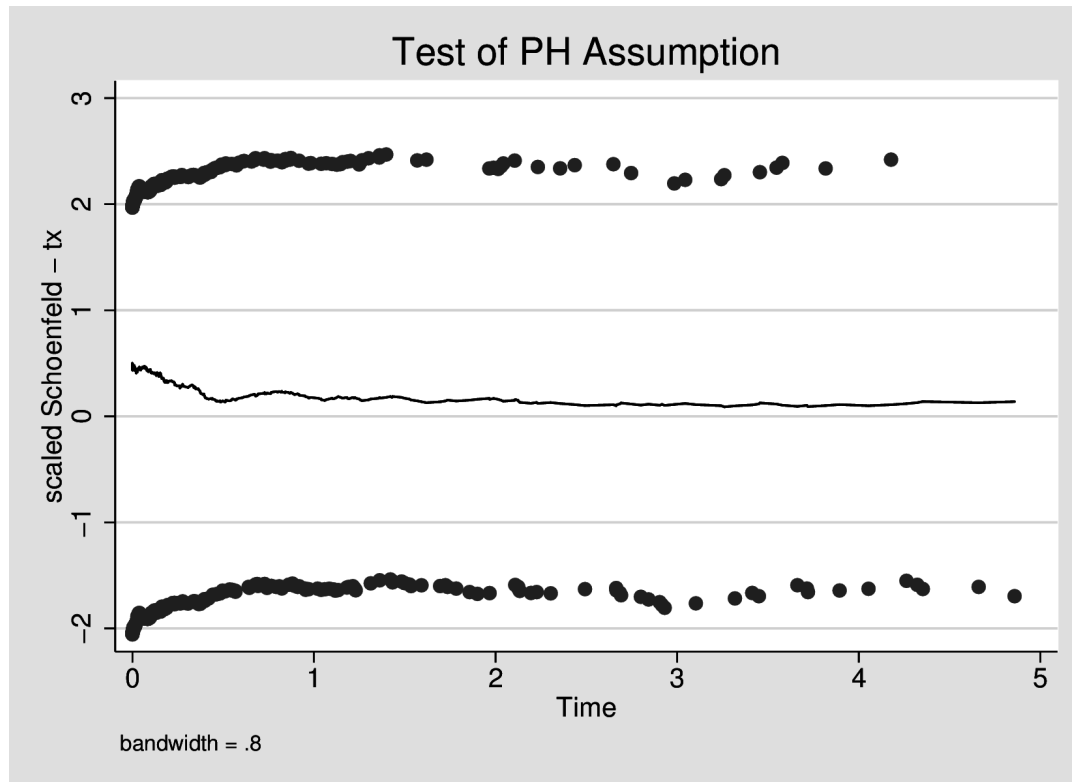
```
. estat phtest, plot(tx) name(phtest, replace) scheme(s2mono)
```

```
. 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 ($p = 0.14$) 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:

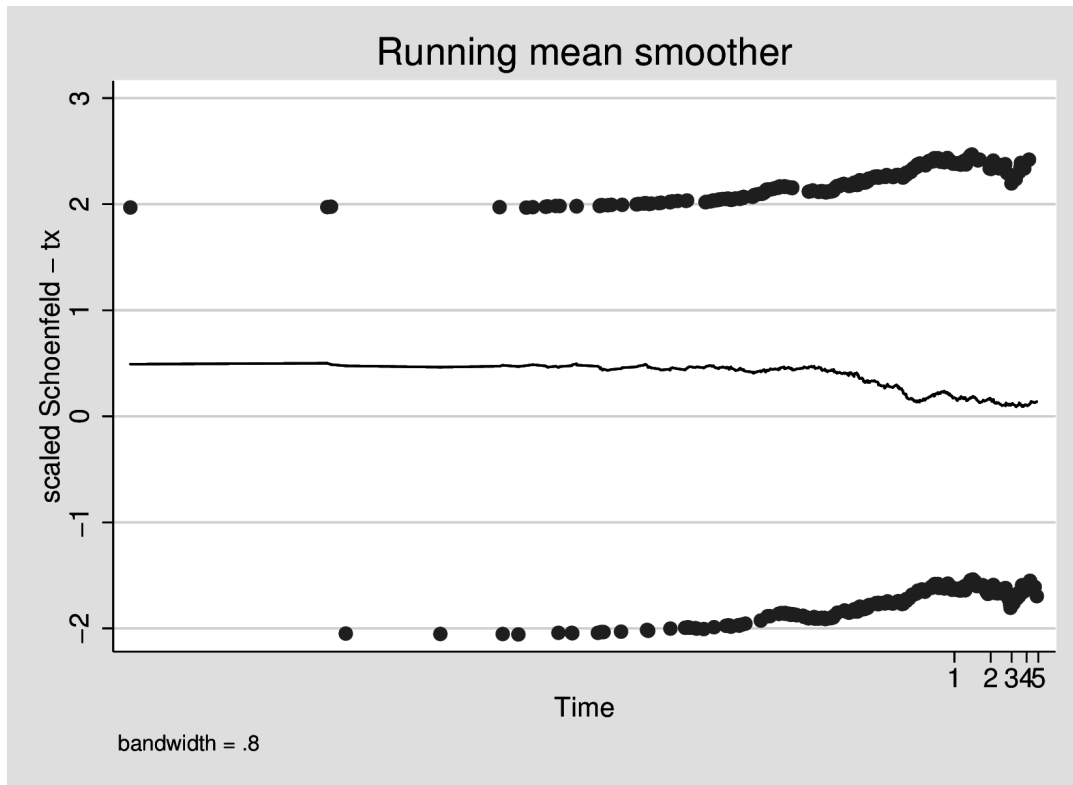
```
. estat phtest, detail log
```

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

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

	rho	chi2	df	Prob>chi2
tx	-0.10840	5.31	1	0.0212
global test		5.31	1	0.0212

```
. 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
```



Again, there is little evidence for non-proportionality ($p = 0.15$).

(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)
(151 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 =          520                Number of obs =          671
No. of failures =          461
Time at risk    = 528.6147698
Log likelihood  = -2527.2463                LR chi2(3) =          15.98
                                                Prob > chi2 =          0.0011
```

	_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
	sex	1.030271	.1027134	0.30	0.765	.8474039 1.2526
	1.tx	1.511383	.1604353	3.89	0.000	1.227491 1.860933
	1.timeband	1	(omitted)			

```
tx#timeband |
      1 1 | .5775587 .1351058 -2.35 0.019 .3651558 .9135115
```

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

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

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

```
No. of subjects =          520          Number of obs =          671
No. of failures =          461
Time at risk    = 528.6147698
Log likelihood  = -2527.2463          LR chi2(3) =          15.98
                                          Prob > chi2 =          0.0011
```

```
-----
      _t | Haz. Ratio   Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
      tx |   1.511383   .1604353    3.89   0.000    1.227491   1.860933
      sex |   1.030271   .1027134    0.30   0.765    .8474039   1.2526
      |
      c.tx#|
      c.timeband | .5775587 .1351058 -2.35 0.019 .3651558 .9135115
-----
```

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

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

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

```
No. of subjects =          520          Number of obs =          671
No. of failures =          461
Time at risk    = 528.6147698
Log likelihood  = -2527.2912          LR chi2(2) =          15.89
                                          Prob > chi2 =          0.0004
```

```
-----
      _t | Haz. Ratio   Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
timeband#c.tx |
      0 |   1.514314   .1604482    3.92   0.000    1.230347   1.863822
      1 |   .8762385   .1826613   -0.63   0.526    .5823448   1.318452
-----
```

This model provides little or no evidence that the hazard ratio is time-dependent ($p = 0.57$). The hazard ratio in the first year is 1.83 (95% CI: 1.48, 2.25), while the hazard ratio after the first year is 1.60 (95% CI: 1.08, 2.39).

(iii)

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

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

```
      failure _d: event
```

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

Cox regression -- no ties

```
No. of subjects =          520          Number of obs =          671
No. of failures =          461
Time at risk   = 528.6147698
Log likelihood = -2527.2912          LR chi2(2)   =          15.89
                                          Prob > chi2 =          0.0004
```

	_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
-----+-----						
main						
	tx	1.514314	.1604482	3.92	0.000	1.230347 1.863822
-----+-----						
tvc						
	tx	.5786372	.1353098	-2.34	0.019	.3658975 .9150677

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

Again, we find little evidence for a time-dependent hazard ratio ($p = 0.57$). 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 =          520          Number of obs =          671
No. of failures =          461
Time at risk   = 528.6147698
Log likelihood = -2526.2791          LR chi2(2)   =          17.91
                                          Prob > chi2 =          0.0001
```

	_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
-----+-----						
main						
	tx	1.609113	.1836904	4.17	0.000	1.28652 2.012596
-----+-----						
tvc						
	tx	.758663	.0783489	-2.67	0.007	.6196455 .928869

Note: variables in tvc equation interacted with _t

The interpretation of this model is as follows: the hazard ratio at time 0 is 1.97 (95% CI: 1.57, 2.48); for each year since diagnosis, the rate tends to decrease by 1-0.84=16% (RR=0.84, 95% CI: 0.68, 1.05), although this trend is not significant ($p = 0.12$, as per the Schoenfeld test).

(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 = -1162.4155          Number of obs   =          671
```

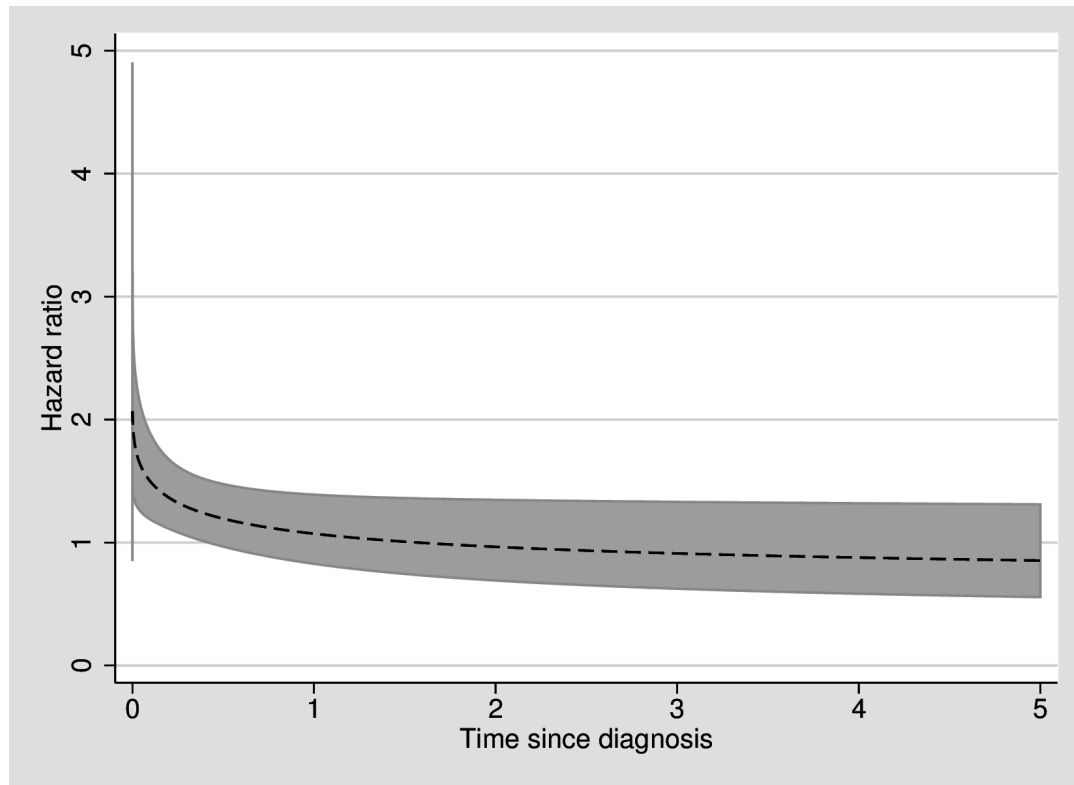
```
-----+-----  
          |      exp(b)   Std. Err.      z    P>|z|      [95% Conf. Interval]  
-----+-----  
xb          |  
    tx      |   1.513752   .1642106     3.82   0.000     1.223819    1.872373  
   _rcs1    |   3.414263   .3157557    13.28   0.000     2.84824    4.092769  
   _rcs2    |   .9859172   .0718682    -0.19   0.846     .8546579    1.137335  
   _rcs3    |   1.027571   .0281718     0.99   0.321     .9738129    1.084298  
   _rcs4    |   1.004915   .0175798     0.28   0.779     .971043    1.039968  
  _rcs_tx1  |   .8664697   .0988198    -1.26   0.209     .6929068    1.083507  
  _rcs_tx2  |   1.078165   .0901261     0.90   0.368     .9152331    1.270103  
   _cons    |   .5161856   .0428399    -7.97   0.000     .4386943    .607365  
-----+-----
```

```
. test _rcs_tx1 _rcs_tx2
```

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

```
      chi2( 2) =      6.43  
      Prob > chi2 =      0.0401
```

```
. 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 limited evidence for time-dependent hazards ($p = 0.37$ 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.