# Biostat III Examination 2016 Answers

## Mark Clements

### May 17, 2016

## Set-up

```
. global folder 5
. 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)      =     547.51
                                                Prob > chi2     =     0.0000
Log likelihood = -888.07465                     Pseudo R2       =     0.2356


------------------------------------------------------------------------------
          lc |        IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
         sex |   2.101776   .1927428     8.10   0.000     1.756011    2.515624
         age |   1.095354   .0045096    22.12   0.000     1.086551    1.104229
       _cons |   1.85e-06   5.02e-07   -48.72   0.000     1.09e-06    3.15e-06
       ln(pt) |          1  (exposure)
```

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

Poisson regression                           Number of obs   =         360
                                             LR chi2(2)      =     1303.83
                                             Prob > chi2     =      0.0000
Log likelihood = -509.91407                  Pseudo R2       =      0.5611


--------------------------------------------------------------------------------
         lc |      IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+-------------------------------------------------------------------
    smoking |  15.63405   1.790159    24.01   0.000     12.49124    19.56759
        age |  1.100899   .0045774    23.12   0.000     1.091964    1.109907
      _cons |  4.70e-07   1.34e-07   -51.21   0.000     2.69e-07    8.21e-07
     ln(pt) |         1  (exposure)
--------------------------------------------------------------------------------

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

Poisson regression                           Number of obs   =         360
                                             LR chi2(2)      =      586.93
                                             Prob > chi2     =      0.0000
Log likelihood = -868.36147                  Pseudo R2       =      0.2526


--------------------------------------------------------------------------------
         lc |      IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+-------------------------------------------------------------------
   asbestos |  3.556892   .3804806    11.86   0.000     2.884149    4.386556
        age |  1.093873   .0044806    21.90   0.000     1.085126     1.10269
      _cons |  2.58e-06   6.76e-07   -49.18   0.000     1.55e-06    4.32e-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)      =      527.17
                                             Prob > chi2     =      0.0000
Log likelihood = -898.24495                  Pseudo R2       =      0.2269


--------------------------------------------------------------------------------
         lc |      IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+-------------------------------------------------------------------
        sex |  2.080201   .1907444     7.99   0.000     1.738017    2.489753
            |
      ageQ5 |
          1 |         1  (base)
          2 |  2.562157   .4319319     5.58   0.000     1.841234    3.565352
          3 |  6.721469   1.053021    12.16   0.000     4.944364    9.137301
          4 |  13.40016   2.123278    16.38   0.000      9.82281    18.28034
```

```
        5 |    20.4727   4.041424    15.29   0.000     13.90412    30.14443
          |
    _cons |   .0000887   .0000132   -62.51   0.000     .0000662    .0001188
    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)      =     564.98
                                                Prob > chi2     =     0.0000
Log likelihood = -879.33904                     Pseudo R2       =     0.2431


-------------------------------------------------------------------------------
       lc |       IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+------------------------------------------------------------------
      sex |   2.085053   .1911968     8.01   0.000     1.742059    2.495578
 ageSpline1 |   1.112402   .0236588     5.01   0.000     1.066985    1.159752
 ageSpline2 |   1.029289   .0634749     0.47   0.640      .912105    1.161528
 ageSpline3 |   .8037657   .1399861    -1.25   0.210     .5713231    1.130778
    _cons |   7.62e-07    7.95e-07   -13.50   0.000     9.84e-08    5.89e-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)      =    1435.18
                                                Prob > chi2     =     0.0000
Log likelihood = -444.23989                     Pseudo R2       =     0.6176


-------------------------------------------------------------------------------
       lc |       IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+------------------------------------------------------------------
      age |   1.103361   .0046225    23.48   0.000     1.094338    1.112459
      sex |   1.475907   .1372634     4.19   0.000     1.229971    1.771019
  smoking |   15.13076   1.739477    23.63   0.000     12.07825    18.95472
 asbestos |   3.443032   .3718914    11.45   0.000     2.786124    4.254825
    _cons |   2.84e-07    8.41e-08   -50.98   0.000     1.59e-07    5.08e-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)

3

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
  smoking  |        0           1 |    Total
-----------+----------------------+----------
         0 | 252.38654   19.281283 | 271.66783
           |     92.90        7.10 |    100.00
-----------+----------------------+----------
         1 | 81.849054   6.4831182 | 88.332172
           |     92.66        7.34 |    100.00
-----------+----------------------+----------
     Total |  334.2356   25.764401 |       360
           |     92.84        7.16 |    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 $\beta_5$, and $I(test)$ 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)      =    1447.96
                                                Prob > chi2     =     0.0000
Log likelihood = -437.84669                     Pseudo R2       =     0.6231
```

4

```
--------------------------------------------------------------------------------
         lc |       IRR   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+-------------------------------------------------------------------
        age |   1.103094   .0046206    23.42   0.000     1.094075    1.112188
        sex |   1.473861   .1368159     4.18   0.000     1.228686    1.767959
  1.smoking |    19.3376   2.728971    20.99   0.000     14.66489    25.49919
 1.asbestos |   7.027045   1.516922     9.03   0.000     4.602826    10.72805
            |
    smoking#|
   asbestos |
        1 1 |   .4007745   .0999591    -3.67   0.000     .2458089     .653435
            |
      _cons |   2.35e-07   7.16e-08   -50.17   0.000     1.30e-07    4.27e-07
      ln(pt) |         1  (exposure)
--------------------------------------------------------------------------------
. est store ModelB
. lrtest ModelA ModelB

Likelihood-ratio test                          LR chi2(1)  =      12.79
(Assumption: ModelA nested in ModelB)          Prob > chi2 =     0.0003
```

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) |   .0082793   .0009574   -41.46   0.000     .0066003     .0103855
--------------------------------------------------------------------------------
```

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
             |     Margin   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
     smoking |
          0  |   .0006102    .0000941     6.48   0.000     .0004257    .0007946
          1  |   .0056096    .0004958    11.31   0.000     .0046378    .0065813
             |
    asbestos |
          0  |   .0015459    .0000953    16.22   0.000     .0013592    .0017327
          1  |   .0046738    .0004901     9.54   0.000     .0037133    .0056344
             |
    smoking#|
    asbestos |
         0 0 |    .000152    .0000209     7.27   0.000      .000111     .000193
         0 1 |   .0010683    .0001859     5.75   0.000      .000704    .0014326
         1 0 |   .0029398    .0001859    15.82   0.000     .0025756    .0033041
         1 1 |   .0082793    .0009574     8.65   0.000     .0064028    .0101558
-------------------------------------------------------------------------------
```

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 = $folder"
Folder = 5
. import delimited "http://biostat3.net/download/exams/2016/$folder/survival.cs
> v", clear
(8 vars, 522 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 |     252    63.37966    .6228296    9.887114    62.15302     64.6063
       1 |     270     62.9975    .5749281    9.447033    61.86557    64.12943
---------+---------------------------------------------------------------------
combined |     522    63.18199    .4225693    9.654575    62.35184    64.01214
---------+---------------------------------------------------------------------
    diff |            .3821607    .8462883               -1.280404    2.044725
-------------------------------------------------------------------------------
    diff = mean(0) - mean(1)                                      t =   0.4516
Ho: diff = 0                                    degrees of freedom =      520

    Ha: diff < 0                 Ha: diff != 0                 Ha: diff > 0
 Pr(T < t) = 0.6741        Pr(|T| > |t|) = 0.6518         Pr(T > t) = 0.3259
. ranksum age, by(tx)
```

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

         tx |      obs    rank sum    expected
------------+-----------------------------------
          0 |      252       66577       65898
          1 |      270       69926       70605
------------+-----------------------------------
   combined |      522      136503      136503

unadjusted variance   2965410.00
adjustment for ties         0.00
                      ----------
adjusted variance     2965410.00

Ho: age(tx==0) = age(tx==1)
           z =    0.394
   Prob > |z| =    0.6934
```

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

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

           |          sex
        tx |        0          1 |     Total
-----------+----------------------+----------
         0 |       92        160 |       252
           |    36.51      63.49 |    100.00
-----------+----------------------+----------
         1 |       92        178 |       270
           |    34.07      65.93 |    100.00
-----------+----------------------+----------
     Total |      184        338 |       522
           |    35.25      64.75 |    100.00

        Pearson chi2(1) =    0.3383   Pr = 0.561
. tab tx smoking, chi row

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

           |        smoking
        tx |        0          1 |     Total
-----------+----------------------+----------
         0 |       47        205 |       252
           |    18.65      81.35 |    100.00
-----------+----------------------+----------
```

```
            1 |          46         224 |        270
              |       17.04       82.96 |     100.00
  -----------+--------------------+----------
        Total |          93         429 |        522
              |       17.82       82.18 |     100.00

          Pearson chi2(1) =    0.2318    Pr = 0.630
. tab tx asbestos, chi row

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

              |        asbestos
         tx |          0           1 |      Total
  -----------+--------------------+----------
          0 |         196          56 |        252
              |       77.78       22.22 |     100.00
  -----------+--------------------+----------
          1 |         215          55 |        270
              |       79.63       20.37 |     100.00
  -----------+--------------------+----------
        Total |         411         111 |        522
              |       78.74       21.26 |     100.00

          Pearson chi2(1) =    0.2670    Pr = 0.605
```

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


--------------------------------------------------------------------------------
     522  total observations
       0  exclusions
--------------------------------------------------------------------------------
     522  observations remaining, representing
     522  subjects
     459  failures in single-failure-per-subject data
 538.4558  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
-------------------------------------------------

        |   Events           Events
tx      |  observed         expected
------+-------------------------
0       |       221           238.29
1       |       238           220.71
------+-------------------------
Total |        459           459.00

          chi2(1) =        2.61
          Pr>chi2 =      0.1059
. 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       83       166         0.3371     0.0299     0.2791    0.3960
        2       56        27         0.2261     0.0266     0.1761    0.2800
        3       39        17         0.1562     0.0232     0.1140    0.2045
        4       29         7         0.1258     0.0213     0.0877    0.1710
        5       24         4         0.1078     0.0201     0.0725    0.1510
tx=1
        1       71       195         0.2721     0.0273     0.2200    0.3267
        2       45        24         0.1768     0.0237     0.1332    0.2257
        3       29        13         0.1239     0.0207     0.0870    0.1678
        4       26         2         0.1151     0.0201     0.0794    0.1580
        5       19         4         0.0963     0.0189     0.0633    0.1373
-------------------------------------------------------------------------------
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 =            522                  Number of obs   =        522
No. of failures =            459
Time at risk    =   538.4557975
                                                 LR chi2(3)      =       3.64
Log likelihood  =    -2529.8439                  Prob > chi2     =     0.3035


------------------------------------------------------------------------------
         _t | Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+-----------------------------------------------------------------
         tx |   1.154286    .1086023     1.53   0.127     .9599027    1.388033
        sex |   .9927266    .0964916    -0.08   0.940     .8205293    1.201061
        age |   .9950297    .0048949    -1.01   0.311     .9854821     1.00467
------------------------------------------------------------------------------
```

10

```
. stcox tx sex, nolog

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

Cox regression -- no ties

No. of subjects =            522              Number of obs   =        522
No. of failures =            459
Time at risk    =   538.4557975
                                              LR chi2(2)      =       2.62
Log likelihood  =    -2530.3543              Prob > chi2     =     0.2704

------------------------------------------------------------------------------
          _t |  Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
          tx |   1.163586    .1090928     1.62   0.106     .968263    1.398309
         sex |   .9935604    .0965671    -0.07   0.947     .8212274   1.202057
------------------------------------------------------------------------------
. stcox tx age, nolog

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

Cox regression -- no ties

No. of subjects =            522              Number of obs   =        522
No. of failures =            459
Time at risk    =   538.4557975
                                              LR chi2(2)      =       3.63
Log likelihood  =    -2529.8467              Prob > chi2     =     0.1627

------------------------------------------------------------------------------
          _t |  Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
          tx |   1.153719    .1082854     1.52   0.128     .9598603   1.386729
         age |   .9950323    .0048953    -1.01   0.311     .9854837   1.004673
------------------------------------------------------------------------------
. stcox tx, nolog

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

Cox regression -- no ties

No. of subjects =            522              Number of obs   =        522
No. of failures =            459
Time at risk    =   538.4557975
                                              LR chi2(1)      =       2.61
Log likelihood  =    -2530.3565              Prob > chi2     =     0.1061

------------------------------------------------------------------------------
          _t |  Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
          tx |   1.163077    .1087764     1.62   0.106     .9682785   1.397066
```

--------------------------------------------------------------------------------

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           |   -0.06370       1.84       1         0.1745
      -------------+--------------------------------------------------
      global test  |                  1.84       1         0.1745
      ----------------------------------------------------------------
. 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
```

Test of PH Assumption

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:
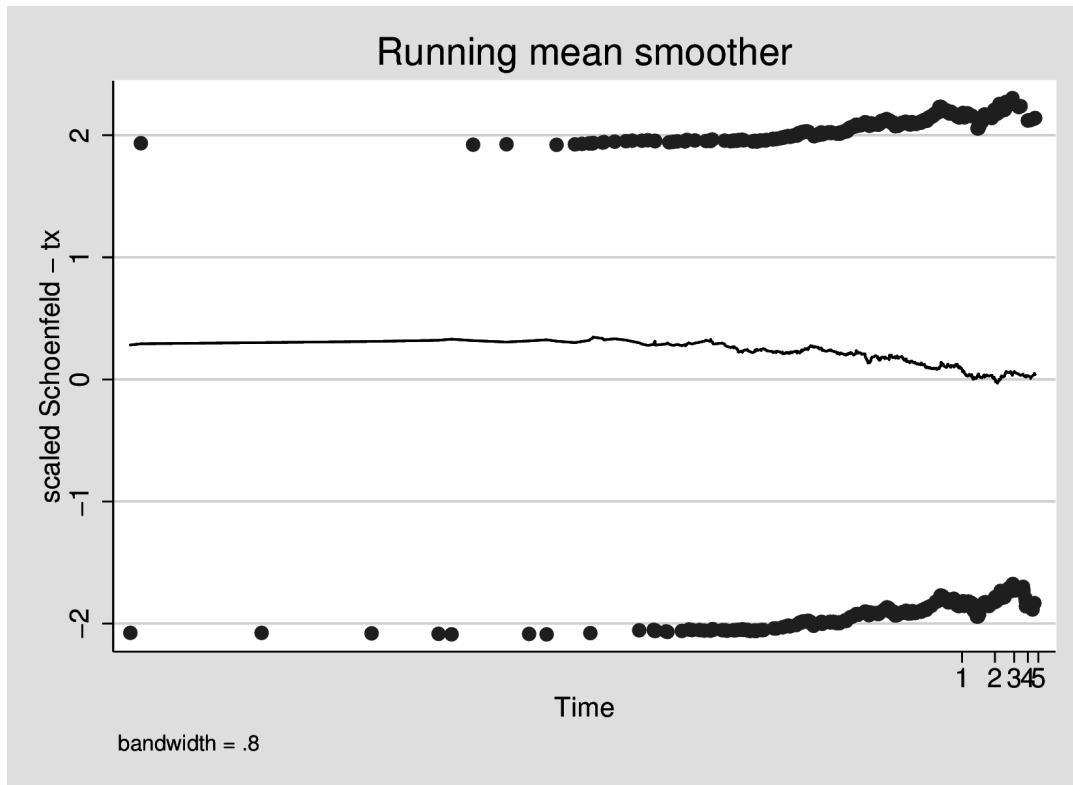
```
. estat phtest, detail log

        Test of proportional-hazards assumption

        Time:  Log(t)
        ----------------------------------------------------------------
                     |       rho        chi2       df      Prob>chi2
        -------------+--------------------------------------------------
        tx           |    -0.02676      0.33        1       0.5684
        -------------+--------------------------------------------------
        global test  |                  0.33        1       0.5684
        ----------------------------------------------------------------
. 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
```

Running mean smoother

bandwidth = .8

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)
(152 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 =          522                 Number of obs   =        674
No. of failures =          459
Time at risk    =   538.4557975
                                               LR chi2(3)      =       4.07
Log likelihood  =    -2529.6271                Prob > chi2     =     0.2540


------------------------------------------------------------------------------
        _t | Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----------+------------------------------------------------------------------
       sex |   .9921265    .0964386    -0.08   0.935     .8200246    1.200348
      1.tx |   1.234292    .1307633     1.99   0.047     1.002859    1.519133
 1.timeband |   20.09028          .       .       .            .           .
           |
```

14

```
tx#timeband |
      1 1  |    .7588576    .1741348    -1.20   0.229     .4839889    1.189831
------------------------------------------------------------------------------
. stcox tx sex c.tx#c.timeband, nolog

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

Cox regression -- no ties

No. of subjects =          522                   Number of obs   =        674
No. of failures =          459
Time at risk    =   538.4557975
                                                 LR chi2(3)      =       4.07
Log likelihood  =   -2529.6271                   Prob > chi2     =     0.2540


------------------------------------------------------------------------------
         _t |  Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+-----------------------------------------------------------------
         tx |   1.234292    .1307633     1.99   0.047     1.002859    1.519133
        sex |   .9921265    .0964386    -0.08   0.935     .8200246    1.200348
            |
       c.tx#|
 c.timeband |   .7588576    .1741348    -1.20   0.229     .4839889    1.189831
------------------------------------------------------------------------------
. stcox c.tx#i.timeband, nolog

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

Cox regression -- no ties

No. of subjects =          522                   Number of obs   =        674
No. of failures =          459
Time at risk    =   538.4557975
                                                 LR chi2(2)      =       4.06
Log likelihood  =   -2529.6304                   Prob > chi2     =     0.1311


------------------------------------------------------------------------------
         _t |  Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
--------------+---------------------------------------------------------------
timeband#c.tx |
          0 |   1.233572    .1303855     1.99   0.047     1.002754     1.51752
          1 |   .9363188    .1906887    -0.32   0.747     .6281595    1.395654
------------------------------------------------------------------------------
```

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 `tvc` and `texp` options:

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

        failure _d:  event
```

```
    analysis time _t:  tsurv
               id:  id


Cox regression -- no ties

No. of subjects =           522                  Number of obs   =          674
No. of failures =           459
Time at risk    =  538.4557975
                                                 LR chi2(2)      =        4.06
Log likelihood  =   -2529.6304                   Prob > chi2     =      0.1311


------------------------------------------------------------------------------
        _t | Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----------+------------------------------------------------------------------
main       |
        tx |   1.233572    .1303855     1.99   0.047     1.002754     1.51752
-----------+------------------------------------------------------------------
tvc        |
        tx |   .7590309    .1741616    -1.20   0.230     .4841156    1.190063
------------------------------------------------------------------------------
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 =           522                  Number of obs   =          674
No. of failures =           459
Time at risk    =  538.4557975
                                                 LR chi2(2)      =        4.47
Log likelihood  =   -2529.4288                   Prob > chi2     =      0.1072


------------------------------------------------------------------------------
        _t | Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----------+------------------------------------------------------------------
main       |
        tx |   1.266262     .14277     2.09   0.036     1.015199    1.579413
-----------+------------------------------------------------------------------
tvc        |
        tx |   .8718413    .0885579    -1.35   0.177     .7144569    1.063895
------------------------------------------------------------------------------
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
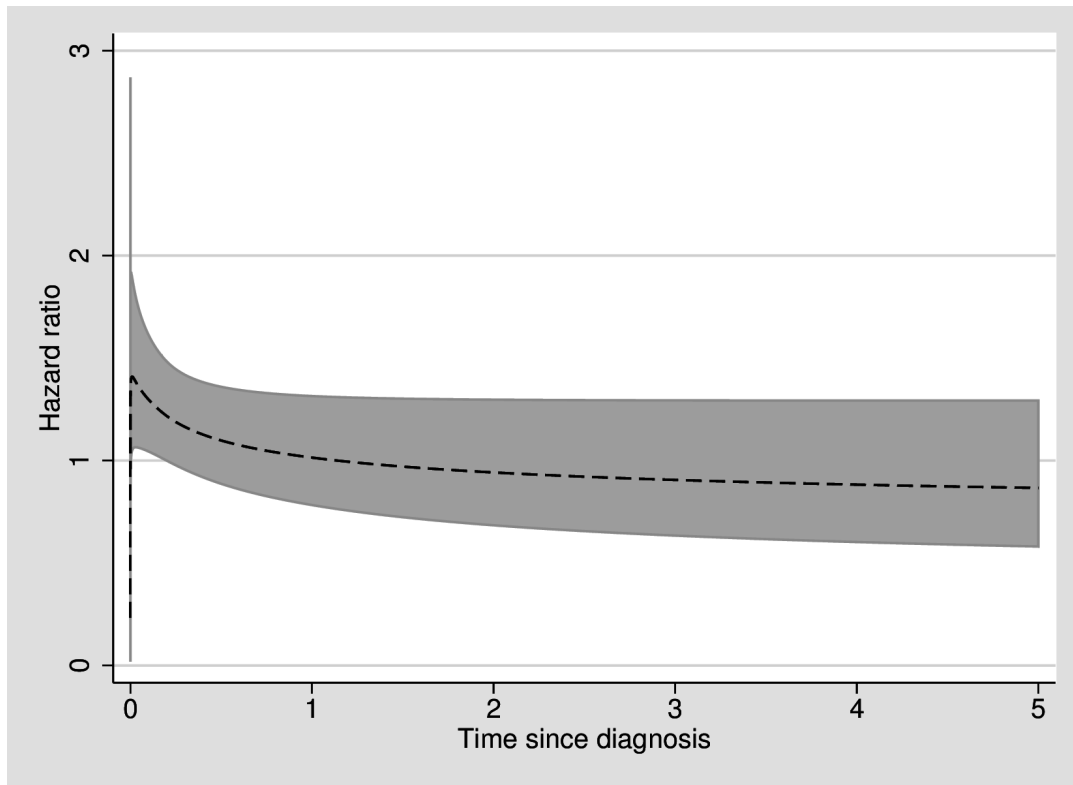
ratio:

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

Log likelihood = -1195.3843                     Number of obs   =        674


------------------------------------------------------------------------------
             |     exp(b)   Std. Err.      z    P>|z|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
xb           |
          tx |   1.184981   .1266655     1.59   0.112     .9610019    1.461162
       _rcs1 |   3.020649   .2299325    14.52   0.000     2.601996    3.506662
       _rcs2 |   .9115485   .0544966    -1.55   0.121     .8107575    1.024869
       _rcs3 |    1.03003   .0282722     1.08   0.281     .9760816     1.08696
       _rcs4 |   .9979048   .0169522    -0.12   0.902      .965226     1.03169
    _rcs_tx1 |   1.074447   .1239423     0.62   0.534     .8570282    1.347023
    _rcs_tx2 |   1.170552   .1053925     1.75   0.080     .9811864    1.396465
       _cons |   .5874006   .0462969    -6.75   0.000     .5033216    .6855249
------------------------------------------------------------------------------
. test _rcs_tx1 _rcs_tx2

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

           chi2(  2) =    3.51
         Prob > chi2 =    0.1732
. 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.