# Survival Analysis (Chapter 7)

- Survival (time-to-event) data

- Kaplan-Meier (KM) estimate/curve

- Log-rank test

- Proportional hazard models (Cox regression)

- Parametric regression models

# Survival Data: Features

- Time-to-event ("event" is not always death)

| | Initiating event | Terminating event |
|---|---|---|
| Neuroblastoma Pt. | Treatment | Death |
| HIV+ Subject | Randomization | HIV-Disease |
| Leukemia Pt. | Enrollment | Relapse |
| Transplant Pt. | Transplant | Death |

- One "event" per person (there are models to handle multiple events per person)
- Follow-up ends with event
- Time-to-death, Time-to-failure, Time-to-event (used interchangeably)

# Survival Data: Structure

For the $i$th sample, we observe:

$T_i$ = time in days/weeks/months/… since origination of the study/treatment/…

$$\delta_i = \begin{cases} 1, & \text{having event at } T_i \\ 0, & \text{no event as of } T_i \end{cases}$$

$X_i$: covariate(s), e.g., treatment, demographic information

Note: in survival analysis, both $T_i$ and $\delta_i$ are outcomes, i.e., $Y_i = (T_i, \delta_i)$.

Censoring: Some lifetimes are known to have occurred only within certain intervals.

Truncation: We only observe subjects whose event time lies within a certain observational window $(T_L, T_R)$. We have no information on subjects whose event time is not in this interval. (For censored data, we have at least partial information on each subject)

# Censoring

Let T = failure time, and C = censoring time

- Right censoring: T > C (a survival time is not known exactly but known to be greater than some value)

    e.g., lost to follow-up, end of study

- Left censoring: the failed subject is never under observation. It is only known that the subject failed between (0, C).

    e.g., study time to employment, some individuals were already employed at the beginning of the study

- Interval censoring: we do not observe exactly when failure occurred, only that it occurred between time $(C_1, C_2)$.

    e.g., longitudinal study with periodic follow-up and the patient's event time is only known to fall in an interval (L , R].

    – Mathematically the same as left censoring.

Note: we assume censored subjects "not different" in risk.

# Truncation

- Left truncation: similar to left censoring, but we don't know those individuals who failed before time C. (often refer to a delayed entry)

  e.g., exposure to some disease, diagnosis of a disease, entry into a retirement home. Any subjects who experience the event of interest prior to the truncation time are not observed.

- Interval truncation: handled similarly as left truncation

- Right truncation: indistinguishable from right censoring (because failure is certain to occur eventually)

# Censoring and Truncation

- SAS:

    model (time_enter time_end)*failed(0) = …

- Stata:

    stset time_end, failure(failed) enter(time time_enter)

# Goals of a Survival Analysis

- Summarize the distribution of survival times

    - Tool: Kaplan-Meier curves

- Compare the survival between groups, e.g., two treatments in clinical trial

    - Tool: Logrank test

- Understand predictors of survival

    - Tool: Cox regression model/parametric models

# Leukemia Example

Treatment for leukemia: Time is measured from remission from induction therapy until relapse. (VGSM Chapter 3.5)

**Table 3.11.** Weeks in Remission Among Leukemia Patients

Placebo: 1,1,2,2,3,4,4,5,5,8,8,8,8,11,11,12,
12,15,17 22,23

6-MP:   6,6,6,6*,7,9*,10,10*,11*,13,16,17*,
19*,20*,22,23,25*,32*,32*,34*,35*

*: The sample was censored.

# Leukemia Example

a. How many of the 6-MP group were censored?

b. What is the longest time until censoring?

c. Is it possible to estimate the MEAN time until relapse in the 6-MP group? How about in the Placebo group?

- What if the first two time points of the Placebo group were censored, would the mean of the Placebo group then be estimable?

- What does this say about the usefulness of the Mean when we have right censored data?

d. What is the likelihood someone relapses AFTER week 3 in the Placebo group or in other words what is the likelihood someone has NOT relapsed before nor including week 3? Same question in the 6-MP group?
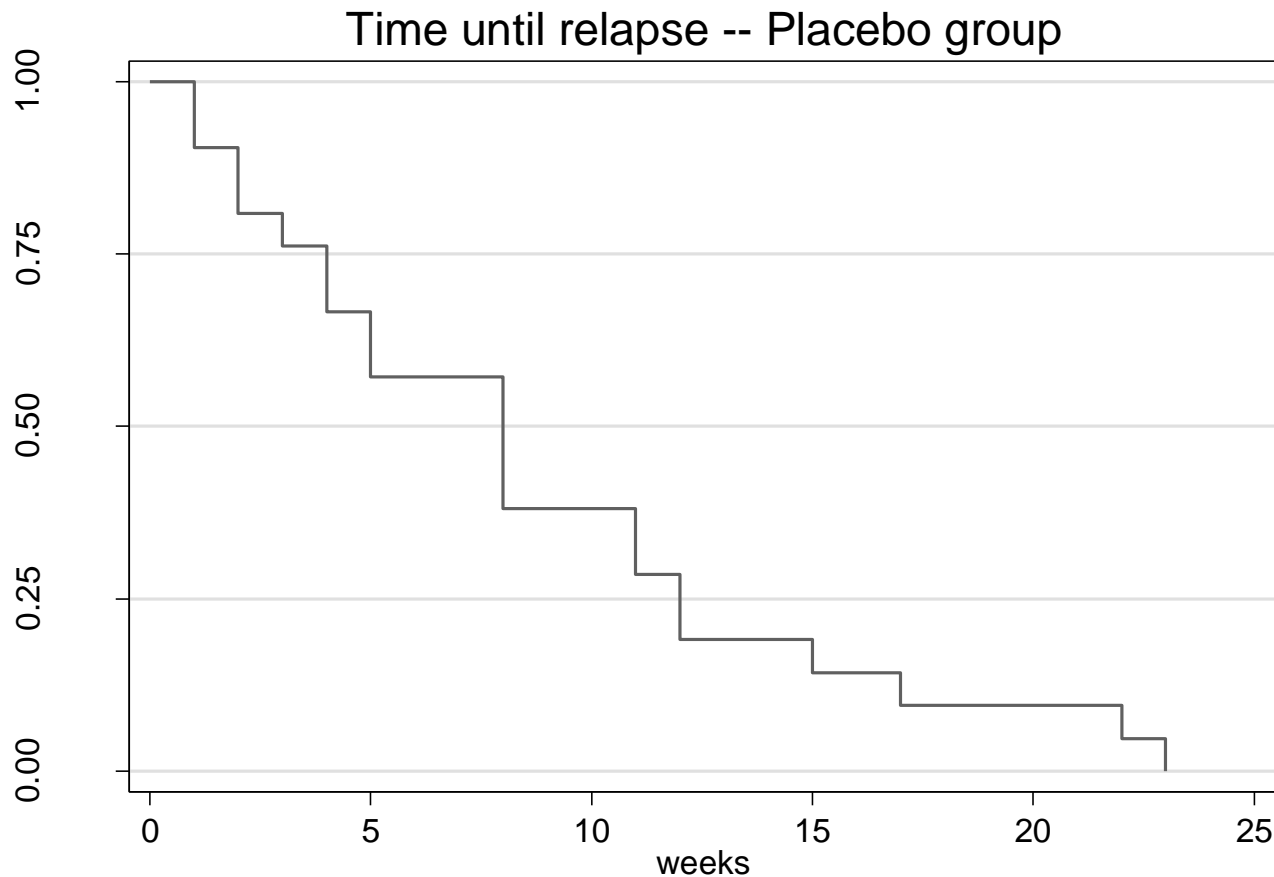
# Survival Function

Definition: The *survival function* at time $t$, denoted $S(t)$, is the probability of being event-free at $t$; equivalently, the probability that the survival time is greater than $t$.

Because there is no censoring in the placebo group, it is simple to estimate the survival probability at each week $t$ by simply taking the percentage of the sample who have not had an event, e.g., S(1)=19/21, S(2)=17/21, ….

In the 6-MP group, because of the right censoring it is not immediately obvious how to estimate the survival probabilities.

- For example, a naïve and mistaken way to estimate the probability of relapse after week 7 (i.e. S(7)) would be to simply consider the person who was censored at week 6 to have instead relapsed at week 6 thus leading to a survival probability of 16/21 = 0.7619, or else to assume that the person censored at week 6 instead has still not relapsed by time 7 and to take 17/21. The first method is too pessimistic and the second is too optimistic.

# Survival Function: without censoring

## Time until relapse -- Placebo group



Use the graph to identify what is median time until relapse in the Placebo group? Would you have come to this same conclusion looking at the raw data?

# Kaplan Meier Estimator

The solution is to rethink the way to estimate the survival probability by noting that the probability can be broking up into the product of probabilities during specific intervals. For any time $t > t_1$,

S(t) = Pr(event occurs after time $t$) = Pr(survive up to time $t_1$)*Pr(survive between time $t_1$ to $t$ | survive up to time $t_1$)

The conditional probability is estimated by using the members of the sample who are still at risk at time $t_1$ (i.e. those who are still known to be at risk at time $t_1$ "the risk set")

# Kaplan Meier Estimator

In general, for $t \in [t_j, t_{j+1}), j = 1, 2, 3, \ldots$, we have:

$$\hat{S}(t) = \left(1 - \frac{d_1}{n_1}\right)\left(1 - \frac{d_2}{n_2}\right) \ldots \left(1 - \frac{d_j}{n_j}\right) = \prod_{i=1}^{j}\left(1 - \frac{d_i}{n_i}\right)$$

where:

$d_i$ = the number of people who have an event during the interval $[t_i, t_{i+1})$

$n_i$ = the number of people at risk just before the beginning of the interval $[t_i, t_{i+1})$

Note that the KM estimator is a step (staircase) function, with the intervals closed at left and open at right.

# Kaplan-Meier Curves

- Useful for exploring survival data

- Plots estimated "survival" probability versus time

- Drops at failure times

- Constant between failures

- Suggests periods of high/low events

- Helps to compare groups

# Survival Function: comparing 3 methods



Time until relapse -- 6MP group

Legend: KM estimate — pessimistic — optimistic

# Survival Function: KM estimate

a. Based on the KM what is the estimated median time until relapse. Would we have been able to find the median using the raw data like in the placebo case?

b. Based on the KM can we estimate at what time there would be 75% of the 6MP group who would have relapsed?

| Time | Beg. Total | Fail | Net Lost | Survivor Function | Std. Error | [95% Conf. Int.] | |
|------|------------|------|----------|-------------------|------------|------------------|--------|
| 6  | 21 | 3 | 1 | 0.8571 | 0.0764 | 0.6197 | 0.9516 |
| 7  | 17 | 1 | 0 | 0.8067 | 0.0869 | 0.5631 | 0.9228 |
| 9  | 16 | 0 | 1 | 0.8067 | 0.0869 | 0.5631 | 0.9228 |
| 10 | 15 | 1 | 1 | 0.7529 | 0.0963 | 0.5032 | 0.8894 |
| 11 | 13 | 0 | 1 | 0.7529 | 0.0963 | 0.5032 | 0.8894 |
| 13 | 12 | 1 | 0 | 0.6902 | 0.1068 | 0.4316 | 0.8491 |
| 16 | 11 | 1 | 0 | 0.6275 | 0.1141 | 0.3675 | 0.8049 |
| 17 | 10 | 0 | 1 | 0.6275 | 0.1141 | 0.3675 | 0.8049 |
| 19 | 9  | 0 | 1 | 0.6275 | 0.1141 | 0.3675 | 0.8049 |
| 20 | 8  | 0 | 1 | 0.6275 | 0.1141 | 0.3675 | 0.8049 |
| 22 | 7  | 1 | 0 | 0.5378 | 0.1282 | 0.2678 | 0.7468 |
| 23 | 6  | 1 | 0 | 0.4482 | 0.1346 | 0.1881 | 0.6801 |
| 25 | 5  | 0 | 1 | 0.4482 | 0.1346 | 0.1881 | 0.6801 |
| 32 | 4  | 0 | 2 | 0.4482 | 0.1346 | 0.1881 | 0.6801 |
| 34 | 2  | 0 | 1 | 0.4482 | 0.1346 | 0.1881 | 0.6801 |
| 35 | 1  | 0 | 1 | 0.4482 | 0.1346 | 0.1881 | 0.6801 |

# Comparing Survival Curves



Time until relapse -- Kaplan-Meier

# Log Rank Test

H0: survival distributions are equal at all followup times.

HA: the two survival curves differ at one or more points in time.

Compares observed number of events in different intervals with expected number assuming two survival curves are the same. (a Chi-square test)

```
Log-rank test for equality of survivor functions


        |   Events          Events
trt     |  observed        expected
------+--------------------------
0       |        21           10.75              chi2(1) =        16.79
1       |         9           19.25              Pr>chi2 =      0.0000
------+--------------------------
Total |         30           30.00
```

Assumptions:
- 2 or more independent groups
- Censoring independent of future risk

Does not assume:
- Proportional hazards (discussed later)
- Equal censoring (censoring process does not depend on group/treatment)

# Log Rank Test: multiple groups ($K > 2$)

- K-group log-rank

- H0: survival curves equal for all groups

- HA: some or all of the survival curves differ at one or more points in time

- Treats K groups as unordered

- Analogous to F-test

- When rejected, unclear interpretation: use KM plots to examine where the important differences arise.

# Hazard function

Definition: The *survival function* at time $t$, denoted $S(t)$, is the probability of being event-free at $t$. The *cumulative incidence function* at time $t$, denoted $F(t) = 1 - S(t)$, is the complementary probability that the event has occurred time by $t$.

Definition: The *hazard function* $h(t)$ is the short-term event rate for subjects who have not yet experienced the outcome event.

- Also known as the instantaneous failure rate, force of mortality, and age-specific failure rate

  $h(t) = -d/dt \log[S(t)] \approx$ #failure at day t /# followed to t

- Similar to conditional probability of failure

- Difficult to calculate: few failures in small time periods (using smoothing technique)

# Example: Pediatric Kidney Transplant

- 9883 kids (age $\leq$ 18) with kidney transplant

- Event: Time from transplant to death

- UNOS database covering 1990-2002

  *38,005 person-years at risk, 465 deaths*

- Does donor source influence survival?

  *cadaveric vs. living donor*
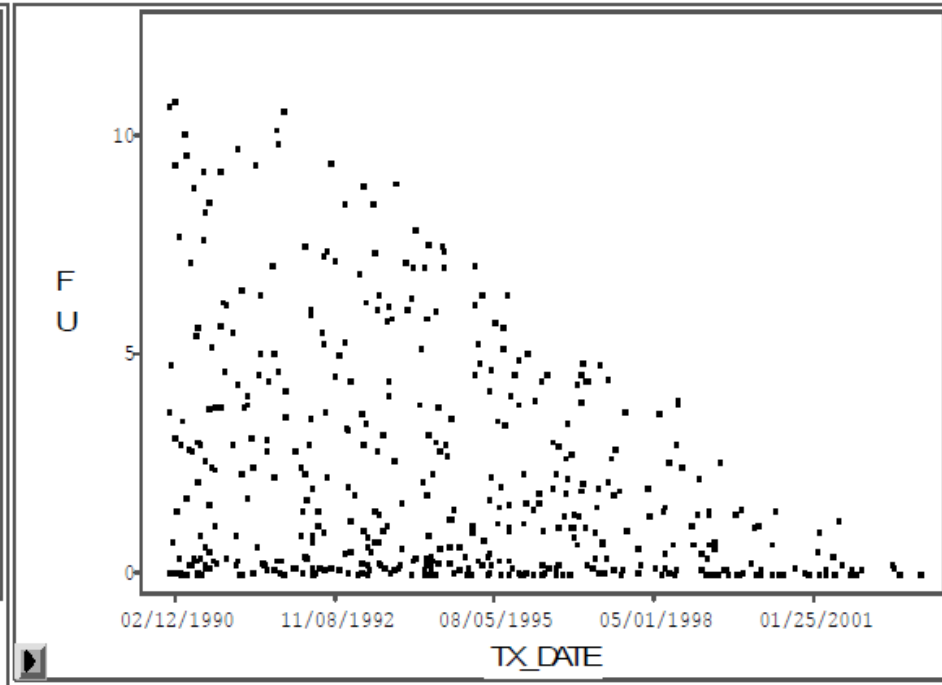
- Does transplant year affect survival?

# Follow-up Time



- Explain why there is a lower triangular shape.
- Explain why there are clumps of observations near the diagonal in the Death = 0 plot

# Survival Curves by Donor Type

## Kaplan-Meier survival estimates



Can you tell where the greatest risk of death is? That is, can you describe the Hazard function?

# Survival Curves by Donor Type

```
Summary Statistics for Time Variable fu
                  Quartile Estimates
            Point            95% Confidence Interval
Percent    Estimate     Transform     [Lower       Upper)
    75        .          LOGLOG          .            .
    50        .          LOGLOG          .            .
    25        .          LOGLOG          .            .


     Mean     Standard Error
   9.9827          0.0493
NOTE: The mean survival time and its standard error were underestimated because the largest
      observation was censored and the estimation was restricted to the largest event time.


Summary of the Number of Censored and Uncensored Values
                                                          Percent
Stratum          txtype      Total   Failed    Censored   Censored
    1                 0        5148      177        4971      96.56
    2                 1        4627      288        4339      93.78
-----------------------------------------------------------------
  Total                        9775      465        9310      95.24


Test of Equality over Strata
                                Pr >
Test       Chi-Square     DF    Chi-Square
Log-Rank     47.0910       1      <.0001
Wilcoxon     46.5318       1      <.0001
-2Log(LR)    50.2749       1      <.0001
```

# Survival Curves by Donor Type

a. Why are there no estimates given under the Quartile Estimates output?

b. Can we conclude that these survival curves are different from one another?

c. Which group has better survival?

# Mortaility Hazard: LOWESS



Fig. 7.1. Mortality Rate for Pediatric Kidney Transplant Recipients

# Hazard for Kidney Transplant Data

- Peaks in first weeks after transplant

- Maximum Hazard: $\approx$ 0.2 deaths/1000 person-days

- Steadily decreases until year 3 (5-fold drop)

- Stabilizes through year 12

- A simple mathematical function of time?

- What does this imply about risk?

# LOWESS-Smoothed Hazard Function by donor type



Fig. 7.2. Smoothed Mortality Rates for Recipients by Kidney Donor Type

# Hazard Ratio

- Relative short-term risk at time $t$: $HR(t) = h_c(t)/h_l(t)$, where:
  $h_c(t)$: hazard function in the recipients of kidneys from recently deceased donors
  $h_l(t)$: hazard function in the recipients of kidneys from living donors

- If $h_c(t) = r*h_l(t)$, **proportional** hazards
  *hazards have same shape*

- Hazards may be complex function of time.

- $r$ can be interpreted as a relative hazard

# Proportional Hazard: kidney data

**Table 7.2.** Smoothed Death Rates (per 1,000 Days) by Donor Type

| Years since transplantation | Smoothed rates | | Death rate ratio |
|---|---|---|---|
| | Cadaveric | Living | |
| 0.25 | 0.235 | 0.098 | 2.40 |
| 0.50 | 0.193 | 0.082 | 2.36 |
| 1.00 | 0.138 | 0.061 | 2.27 |
| 2.00 | 0.088 | 0.038 | 2.30 |
| 3.00 | 0.061 | 0.027 | 2.25 |
| 4.00 | 0.063 | 0.026 | 2.37 |
| 5.00 | 0.065 | 0.032 | 2.03 |

The hazard function for Cadaveric kidney is approximately proportional to the hazard function for living kidney donor.

# Proportional Hazards Assumption

Definition: Under the *proportional hazards assumption* the hazard ratio does not vary with time. That is, $HR(t) \equiv HR$.

Analog: multiple linear regression without interaction terms, e.g., similar age effects in males and females (what is the moderator in the interaction?)

# Cox Proportional Hazards Model

The Cox Proportional Hazards model (CPH) (or "Cox model" or "Cox regression") is the most commonly applied model in medical time-to-event studies. The original reference is: Cox DR. Regression models and life tables (with discussion). *Journal of the Royal Statistical Society Series B* 1972;34:187–220.

The Cox proportional hazards model does not make any assumption about the shape of the underlying hazards, but makes the assumption that the hazards for patient subgroups are proportional over follow-up time.

There are parametric hazard models that assume the hazard function follows a particular functional form derived from, e.g. an exponential or Weibull distribution.
* These parametric models are more efficient if the functional form of the hazard is correct and also allow prediction beyond the data (although risky like extrapolation in general because the form of the hazard can't be verified).

The CPH is much more common because of its robustness to the form of the hazard and because it has been show to be relatively efficient.

# Developing the CPH Regression Model

1. We are interested in modeling the hazard function, $h(t; \mathbf{X})$, for individuals with covariate vector $\mathbf{X}$, where $\mathbf{X}$ represents $k$ predictors $X_1, \ldots, X_k$. (e.g. liver type, age, year of transplant, etc.)

2. The hazard function should be non-negative, so since $\exp(\beta_0 + \beta_1 X_1 + \ldots + \beta_k X_k)$ is always positive, it is natural to model:

$$\log h(t; \mathbf{X}) = \beta_0 + \beta_1 X_1 + \ldots + \beta_k X_k$$

As is, the right hand side of this equation does not depend on $t$ and although in some cases it might be ok to assume the hazard is constant across time, constant hazard usually is unrealistic.

# Developing the CPH Regression Model

3.  The solution is to replace the intercept with a function of time, called the baseline hazard function $h_0(t)$ which is non-parametrically estimated in the Cox model (similar to the non-parametric way the KM estimates the Survival function) and serves as the reference for how the hazard changes over time. Thus the model becomes:

$$\log h(t; \mathbf{X}) = \log h_0(t) + \beta_1 X_1 + \ldots + \beta_k X_k$$

$$h(t; \mathbf{X}) = h_0(t) \exp(\beta_1 X_1 + \ldots + \beta_k X_k)$$

The hazard at any particular time for some covariate combination is proportional to the baseline hazard. Hence comparisons between two sets of covariate values will NOT depend on time. Show that the hazard ratio for a 1 unit increase in X, i.e. h(t;X+1)/h(t;X), does not depend on t.

4.  Estimates for $\beta_1, \ldots, \beta_k$ are log hazard ratios with $\exp(\beta_k)$ representing the hazard ratio for a one unit increase in $X_k$. Note the HR does not depend on $t$ (time). The proportionality assumption (i.e. that the HR does not depend on time) should be checked and we'll see how to do it later.

# Partial Likelihoods

Assuming the event and censoring time are independent and no ties between the event time,

$$P[\text{individual dies at } t_i \mid \text{one death at } t_i]$$

$$= \frac{P[\text{individual dies at } t_i \mid \text{survival to } t_i]}{P[\text{one death at } t_i \mid \text{survival to } t_i]}$$

$$= \frac{b[t_i \mid \mathbf{Z}_{(i)}]}{\sum_{j \in R(t_i)} b[t_i \mid \mathbf{Z}_j]} = \frac{b_0(t_i) \exp[\boldsymbol{\beta}' \mathbf{Z}_{(i)}]}{\sum_{j \in R(t_i)} b_0(t_i) \exp[\boldsymbol{\beta}' \mathbf{Z}_j]}$$

$$= \frac{\exp[\boldsymbol{\beta}' \mathbf{Z}_{(i)}]}{\sum_{j \in R(t_i)} \exp[\boldsymbol{\beta}' \mathbf{Z}_j]}.$$

where $R(t_i)$ is the risk set at time $t_i$, i.e., individuals who are still under study at a time just prior to $t_i$. The partial likelihood is:

$$L(\boldsymbol{\beta}) = \prod_{i=1}^{D} \frac{\exp\left[\sum_{k=1}^{p} \beta_k Z_{(i)k}\right]}{\sum_{j \in R(t_i)} \exp\left[\sum_{k=1}^{p} \beta_k Z_{jk}\right]}$$

Note that this is the same LL as for conditional logistic regression.

When multiple events occurred at the same time, there are multiple ways to calculate or approximate the partile likelihood function. In SAS, use `/ties=`$_\top$; in Stata, specify one of the `efron/breslow/exactm/exactp` options.

# CPH Regression Model: Kidney Data

From the survival curves, it looks like kidneys from live donors lead to better survival than cadaver kidneys. But this is not a randomized experiment so the choice of whether someone got a living or cadaver kidney was NOT randomized. Hence there may be some confounders that need to be controlled when considering which type of kidney is better. For example, there may be some other variables that are predictive of whether someone gets a live kidney and that same variable may also effect survival.

Here we consider how the Age of the recipient and year of transplantation are related to survival.

# Survival Curves by Age Group



Kaplan-Meier survival estimates, by agecat

Legend:
- agecat = <3
- agecat = 3-4
- agecat = 5-6
- agecat = 7-16
- agecat = 17
- agecat = 18

Which age group has the worst prognosis examining the KM curves?

# Survival Curves by Age Group: log-rank test

```
Log-rank test for equality of survivor functions

          |   Events          Events
agecat    |  observed        expected
--------+-----------------------------
<3        |        82           40.37
3-4       |        37           27.35
5-6       |        24           29.42
7-16      |       231          268.36
17        |        46           47.18
18        |        40           47.31
--------+-----------------------------
Total     |       460          460.00

          chi2(5)  =      53.78
          Pr>chi2  =     0.0000
```

# Age Group vs Donor Type

Here is a plot of the age of the recipient versus the proportion who received a kidney from a cadaver.



There is a trend such that the older the recipient, the more likely to receive a cadaver kidney.

If we would "control" for age, how would you expect it to change the effect found for kidney type?

# Survival Curves by Transplant Year



Kaplan-Meier survival estimates, by yearsperf

yearsperf = 1990-1994
yearsperf = 1995-1999
yearsperf = 2000-2002

analysis time

Which group has the worst prognosis examining the KM curves?

# Survival Curves by Transplant Year: log-rank test

```
Summary of the Number of Censored and Uncensored Values

                                                      Percent
Stratum     yearsperf      Total   Failed    Censored   Censored
      1     1990-1994       3651      261        3390      92.85
      2     1995-1999       3929      171        3758      95.65
      3     2000-2002       2195       33        2162      98.50
-----------------------------------------------------------------
   Total                    9775      465        9310      95.24


                              Trend Tests
                 Test       Standard
Test          Statistic       Error      z-Score   Pr > |z|    Pr < z     Pr > z
Log-Rank       -30.2577      13.1027     -2.3093     0.0209     0.0105     0.9895
Wilcoxon     -277528.00    105598.540    -2.6281     0.0086     0.0043     0.9957
```

The trend tests test for trend across the levels of the strata (three time period).
The p-values suggest that there have been secular improvements in the
prognosis for survival after kidney transplant as time passes.

# Transplant Year vs Donor Type

```
           |        txtype
yearsperf  |   Living  Cadaveric |     Total
-----------+---------------------+----------
1990-1994  |    1,808      1,843 |     3,651
           |    49.52      50.48 |    100.00
-----------+---------------------+----------
1995-1999  |    2,099      1,830 |     3,929
           |    53.42      46.58 |    100.00
-----------+---------------------+----------
2000-2002  |    1,241        954 |     2,195
           |    56.54      43.46 |    100.00
-----------+---------------------+----------
    Total  |    5,148      4,627 |     9,775
           |    52.66      47.34 |    100.00

         Pearson chi2(2) =  28.5907   Pr = 0.000
```

We can see in the data that the % of kidneys from cadavers has been going down over time as well. So, the year in which the surgery was performed may also be a confounder for the true effect of kidney type.

Indeed we could even investigate kidney type as a mediator for the improvement in survival over time.

# Cox Proportional Hazards Model: SAS

```
proc phreg data=unos_c;
    model fu*death(0) = txtype/risklimits alpha = .05;
run;
```

| | |
|---|---|
| Data Set | WORK.UNOS_C |
| Dependent Variable | fu |
| Censoring Variable | death |
| Censoring Value(s) | 0 |
| Ties Handling | BRESLOW |

| | |
|---|---|
| Number of Observations Read | 9775 |
| Number of Observations Used | 9775 |

Summary of the Number of Event and Censored Values

| Total | Event | Censored | Percent Censored |
|---|---|---|---|
| 9775 | 465 | 9310 | 95.24 |

Convergence Status

Convergence criterion (GCONV=1E-8) satisfied.

# Cox Proportional Hazards Model: SAS

```
Model Fit Statistics
                    Without           With
Criterion         Covariates        Covariates
-2 LOG L           8023.065          7975.964
AIC                8023.065          7977.964
SBC                8023.065          7982.106


            Testing Global Null Hypothesis: BETA=0
Test                   Chi-Square        DF      Pr > ChiSq
Likelihood Ratio        47.1012          1         <.0001
Score                   47.0867          1         <.0001
Wald                    45.4954          1         <.0001
```

```
                        Analysis of Maximum Likelihood Estimates
                    Parameter      Standard                                  Hazard      95% Hazard Ratio
Parameter    DF     Estimate        Error     Chi-Square   Pr > ChiSq        Ratio     Confidence Limits
txtype       1      0.64466        0.09558      45.4954       <.0001          1.905      1.580      2.298
```

What is the hazard ratio for death in the cadaver group as compared to the living kidney donor group?

# Cox Proportional Hazards Model: Stata

```
. stset fu, failure(death)


      failure event:  death != 0 & death < .
obs. time interval:  (0, fu]
 exit on or before:  failure


--------------------------------------------------------------------------
      9775  total obs.
        25  obs. end on or before enter()
--------------------------------------------------------------------------
      9750  obs. remaining, representing
       461  failures in single record/single failure data
  38004.91  total analysis time at risk, at risk from t =         0
                                earliest observed entry t =         0
                                   last observed exit t =   12.53151
```

# Cox Proportional Hazards Model: Stata

```
. stcox txtype, nohr

Cox regression -- Breslow method for ties

No. of subjects =          9750                    Number of obs   =        9750
No. of failures =           461
Time at risk    =    38004.90961
                                                   LR chi2(1)      =       44.82
Log likelihood  =    -3952.3735                    Prob > chi2     =      0.0000
------------------------------------------------------------------------------
         _t |      Coef.    Std. Err.       z     P>|z|      [95% Conf. Interval]
------------+-----------------------------------------------------------------
     txtype |   .6310981    .0958317      6.59    0.000      .4432714    .8189248
------------------------------------------------------------------------------
```

Note the estimate is slightly different from SAS output.  It is because Stata excludes samples with survival time 0 (fu=0). We can get back these samples by adding a small amount of time.

```
. replace fu=.002 if fu==0
. stset fu, failure(death)        <-- Take a look at Stata output, compare to the previous one
. stcox txtype, nohr
------------------------------------------------------------------------------
         _t |      Coef.    Std. Err.       z     P>|z|      [95% Conf. Interval]
------------+-----------------------------------------------------------------
     txtype |   .6446623    .0955754      6.75    0.000      .457338    .8319866
------------------------------------------------------------------------------
```

# Cox Proportional Hazards Model: Stata

The partial likelihood in CPH model depends on the risk groups at each failure.

- When there is no time-varying covariate, only the relative order of failures matters

- We can shift the origin without changing the estimation results

```
. replace fu=fu+20
. stset fu, failure(death)
. stcox txtype, nohr

------------------------------------------------------------------------------
       _t |      Coef.   Std. Err.      z     P>|z|     [95% Conf. Interval]
----------+-------------------------------------------------------------------
   txtype |   .6446623   .0955754     6.75   0.000      .457338     .8319866
------------------------------------------------------------------------------
```

# Cox Proportional Hazards Model: confounders

```
proc phreg data=unos_c;
    model fu*death(0) = txtype age/risklimits alpha = .05;
run;
```

### Analysis of Maximum Likelihood Estimates

| Parameter | DF | Parameter Estimate | Standard Error | Chi-Square | Pr > ChiSq | Hazard Ratio | 95% Hazard Ratio Confidence Limits | |
|-----------|----|--------------------|----------------|------------|------------|--------------|------------|------|
| txtype | 1 | 0.69376 | 0.09616 | 52.0543 | <.0001 | 2.001 | 1.657 | 2.416 |
| age | 1 | -0.04908 | 0.00854 | 33.0671 | <.0001 | 0.952 | 0.936 | 0.968 |

What is the hazard ratio for death in the cadaver group as compared to the living kidney donor group? Compare to the txtype only model.

```
proc phreg data=unos_c;
    class yearsperf (ref = "1990-1994");
    model fu*death(0) = txtype yearsperf/risklimits alpha = .05;
run;
```

| Parameter | | DF | Parameter Estimate | Standard Error | Chi-Square | Pr > ChiSq | Hazard Ratio | 95% Hazard Ratio Confidence Limits | |
|-----------|---|----|--------------------|----------------|------------|------------|--------------|------------|------|
| txtype | | 1 | 0.63847 | 0.09563 | 44.5758 | <.0001 | 1.894 | 1.570 | 2.284 |
| yearsperf | 1995-1999 | 1 | -0.12912 | 0.10329 | 1.5627 | 0.2113 | 0.879 | 0.718 | 1.076 |
| yearsperf | 2000-2002 | 1 | -0.38078 | 0.19336 | 3.8780 | 0.0489 | 0.683 | 0.468 | 0.998 |

What is the hazard ratio for death in the cadaver group as compared to the living kidney donor group? Compare to the txtype only model.

# Cox Proportional Hazards Model: age + yearsperf

```
. stcox txtype age i.yearsperf

Cox regression -- Breslow method for ties

No. of subjects =          9766                     Number of obs   =        9766
No. of failures =           464
Time at risk    =   233306.5767
                                                    LR chi2(4)      =       83.45
Log likelihood  =   -3960.7595                      Prob > chi2     =      0.0000

------------------------------------------------------------------------------
         _t | Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+-----------------------------------------------------------------
     txtype |   1.989285    .1913899     7.15   0.000     1.647413    2.402102
        age |   .9521617    .0081282    -5.74   0.000     .9363633    .9682267
            |
   yearsperf |
       1995 |   .8750921    .0905664    -1.29   0.197     .7144303    1.071884
       2000 |   .6833657     .132183    -1.97   0.049     .4677413    .9983908
------------------------------------------------------------------------------
```

What is the hazard ratio for 2000-2002 as compared to 1990-1994 (the reference)? What is the hazard ratio for age?

How did the hazard ratio for kidney type change depending on what was controlled for?

# Adjusted Survival Functions

Similar to Adjusted means (LSmeans) which provide a way to present results from a regression back on the original scale, we can obtain the Adjusted (regression fitted) Survival curves for different combinations of predictors.

# Categorical variable: categories with no events

```
. replace death = 0 if yearsperf==1990
(261 real changes made)

. stset fu, failure(death)

     failure event:  death != 0 & death < .
obs. time interval:  (0, fu]
 exit on or before:  failure

-----------------------------------------------------------------
     9775  total obs.
        0  exclusions
-----------------------------------------------------------------
     9775  obs. remaining, representing
      204  failures in single record/single failure data
 38004.96  total analysis time at risk, at risk from t =        0
                            earliest observed entry t =        0
                                 last observed exit t =  12.53151
```

# Categorical variable: categories with no events

```
. stcox txtype i.yearsperf
------------------------------------------------------------------------------
        _t | Haz. Ratio   Std. Err.      z     P>|z|     [95% Conf. Interval]
-----------+------------------------------------------------------------------
    txtype |   1.953708    .2815818     4.65   0.000     1.472918    2.591437
           |
  yearsperf |
      1995 |   1.57e+10    3.17e+09   116.76   0.000     1.06e+10    2.33e+10
      2000 |   1.31e+10          .        .       .            .           .
------------------------------------------------------------------------------

. stcox txtype ib(2000).yearsperf
------------------------------------------------------------------------------
        _t | Haz. Ratio   Std. Err.      z     P>|z|     [95% Conf. Interval]
-----------+------------------------------------------------------------------
    txtype |   1.953708    .2815818     4.65   0.000     1.472918    2.591437
           |
  yearsperf |
      1990 |   3.90e-22          .        .       .            .           .
      1995 |   1.200451    .2413932     0.91   0.364     .8094316    1.780363
------------------------------------------------------------------------------
```

Avoid using category with no events as the reference group (recall the separation problem in logistic regression. difference here?)

# Assessing the Proportional Hazards Assumption

- The proportional hazards assumption is a strong assumption and its appropriateness should always be assessed.
- The model assumes that the ratio of the hazard functions for any two subgroups (i.e. two groups with different values of the explanatory variable X) is constant over follow-up time.
- Note that it is the hazard ratio which is assumed to be constant. The hazard can vary freely with time (baseline hazard is a function of $t$, $h_0(t)$).
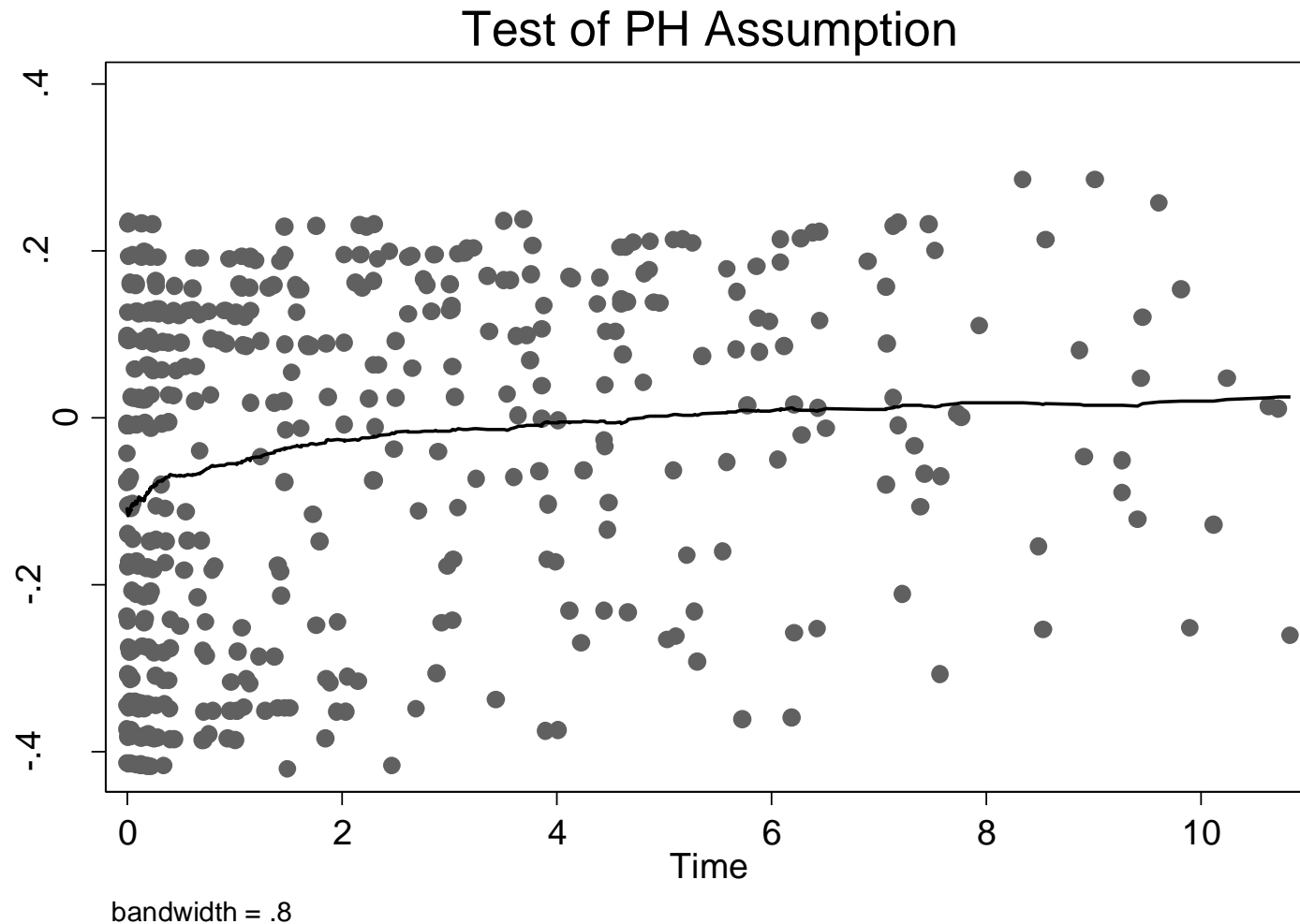
```
. estat phtest, detail
     Test of proportional-hazards assumption
     Time:  Time
     --------------------------------------------------------------
                 |       rho          chi2        df      Prob>chi2
     ------------+-------------------------------------------------
     txtype      |     -0.09261        4.09        1        0.0430
     age         |      0.24686       35.69        1        0.0000
     1990b.year~f|          .             .        1            .
     1995.years~f|      0.05677        1.54        1        0.2142
     2000.years~f|      0.00960        0.04        1        0.8350
     ------------+-------------------------------------------------
     global test |                    39.40        4        0.0000
     --------------------------------------------------------------
```

Can also test the PH assumption on transformation of analysis time (e.g., log or rank).
More on: http://www.ats.ucla.edu/stat/examples/asa/test_proportionality.htm

# Assessing the Proportional Hazards Assumption – using Schoenfeld residuals

```
. estat phtest, plot(age)
```



Test of PH Assumption

bandwidth = .8

# Assessing the Proportional Hazards Assumption – Survival curves

```
.  stcoxkm, by(txtype)
```

# Assessing the Proportional Hazards Assumption – log-log plot

```
. stphplot, by(txtype)
```

# Proportional Hazards Assumption: alternatives

When the proportional hazards assumption is violated, some possible approaches:

- Consider time-varying covariates: X enter the model as a function of time

- Stratification: instead of treating X as a covariate, model the hazard function in each stratum of X

$$h(t; \text{stratum}=j, \mathbf{X}) = h_{0j}(t) \exp(\beta_1 X_1 + \ldots + \beta_k X_k)$$

  We assume that the effect of each of the covariates is the same across strata, but the baseline hazard are different.

- Parametric regression models

# Stratified Analysis

```
. stcox txtype i.yearsperf, strata(agecat)

Stratified Cox regr. -- Breslow method for ties

No. of subjects =            9775                    Number of obs    =       9766
No. of failures =             465
Time at risk    =    38004.95961
                                                    LR chi2(3)       =      62.11
Log likelihood  =    -3272.124                      Prob > chi2      =     0.0000


-----------------------------------------------------------------------------
        _t | Haz. Ratio   Std. Err.       z    P>|z|      [95% Conf. Interval]
-----------+-----------------------------------------------------------------
    txtype |   2.032869    .1962781     7.35   0.000       1.68238    2.456376
           |
 yearsperf |
      1995 |   .8816313    .0912277    -1.22   0.223      .7197937    1.079856
      2000 |   .6717791    .1299931    -2.06   0.040       .459742    .9816096
-----------------------------------------------------------------------------
                                                             Stratified by agecat
```

# Time-depenendent (time-varying) Covariates

Definition: A *time-dependent covariate* in a Cox model is a predictor whose values may vary with time.

- Discrete time-varying covariates:
  - Example: Aurora et al. (1999) followed 124 patients to study the effect of lung transplantation on survival in children with cystic fibrosis. The natural time origin in this study is the time of listing for transplantation, not transplantation itself.
  - Transplantation is then treated as a time-dependent variable.

$$h(t|x) = h_0(t) \exp\{\beta X(t)\}$$
$$= \begin{cases} h_0(t) & \text{before transplantation} \\ h_0(t) \exp(\beta) & \text{at or after transplantation.} \end{cases}$$

  - The data is in the long format, similar to that of longitudinal data.
  - We do not need to change the command syntax: just enter transplantation as a predictor in the model statement.

# Time-depenendent (time-varying) Covariates

- Continuous time-varying covariates: interaction between predictor(s) and function of time, i.e., $z_i(t) = z_i*f(t)$

```
. stcox txtype age i.yearsperf, tvc(txtype age)
                                               LR chi2(6)        =      121.83
Log likelihood  =   -3941.5661                 Prob > chi2       =      0.0000
------------------------------------------------------------------------------
         _t | Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]
------------+-----------------------------------------------------------------
main        |
     txtype |   2.366621    .307691      6.63   0.000     1.834263    3.053485
        age |   .9108782    .0102289    -8.31   0.000     .8910488    .9311488
            |
  yearsperf |
       1995 |   .8758925    .090591     -1.28   0.200      .715177    1.072724
       2000 |   .6845634    .1324103    -1.96   0.050     .4685671    1.000128
------------+-----------------------------------------------------------------
tvc         |
     txtype |   .9287556    .0338322    -2.03   0.042     .8647575    .9974901
        age |   1.020512    .0035547     5.83   0.000     1.013569    1.027503
------------------------------------------------------------------------------
Note: variables in tvc equation interacted with _t
```

*. In SAS, the interaction(s) need to be generated within –proc phreg-:

```
proc phreg data = unos_c;
      class yearsperf (ref = "1990-1994");
      model fu*death(0) = txtype age yearsperf aget txt /risklimits alpha = .05;
      aget = age*fu;
      txt = txtype*fu;
run;
```

# Residual analysis

- Cox-Snell residuals: assessing overall model fit

- Martingale residuals: determining the functional form of covariates

- Deviance residuals: examining model accuracy and identifying outliers

- Schoenfeld/scaled Schoenfeld residuals: checking PH assumption

- Efficient score residuals (Likelihood displacement values, LMAX values, and DFBETAs): identifying influential subjects

# Cox Proportional Hazard Model – Partial Likelihood

An extreme example,

| id | T0 | T | Gender | Failure |
|----|----|----|--------|---------|
| 1  | 0  | 2  | 0      | 1       |
| 2  | 3  | 5  | 1      | 0       |
| 3  | 6  | 8  | 0      | 1       |
| 4  | 9  | 10 | 1      | 1       |

We cannot obtain estimate for gender effect, since there is only one subject at each failure time.  We also cannot obtain KM estimate of the overall survivor function.

Parametric model has to be used with these data.

# Parametric Regression Model

Use a linear model to model log survival time:

$$\log(x) = \mu + \gamma'\mathbf{Z} + \sigma W$$

where W is the error distribution with mean 0 and variance 1.

Choice of distribution:
Exponential
Weibull
log-normal
log-logistic
Gamma
inverse Gaussian
etc.

# Parametric Regression Model

Another representation: accelerated failure-time (AFT) model

$$S(x; Z) = S_0(\exp\{\mathbf{Z}\theta\}x)$$

where $\exp\{\mathbf{Z}\theta\}$ is called an acceleration factor to model the change of time scale compared to the baseline time scale. For example, let Z be the donor type, then the baseline survival function (Z = 0, living donor) is $S_0(x)$, while for cadaveric donor (Z = 1), $S(x; Z=1) = S_0(\exp\{\theta\}x)$. The survival time is accelerated by $e^\theta$ times. Therefore, a positive estimate for $\theta$ implies a shorter survival time.

When $S_0(x) = \exp(\mu + \sigma W)$, the two representations are equivalent with

$$\theta = -\gamma$$

$$
\begin{aligned}
S(x|\mathbf{Z}) &= Pr[X > x|\mathbf{Z}] = Pr[Y > \ln x|\mathbf{Z}] \\
&= Pr[\mu + \sigma W > \ln x - \boldsymbol{\gamma}'\mathbf{Z}|\mathbf{Z}] \\
&= Pr[e^{\mu+\sigma W} > x\exp\{-\boldsymbol{\gamma}'\mathbf{Z}\}|\mathbf{Z}] \\
&= S_0(x\exp\{-\boldsymbol{\gamma}'\mathbf{Z}\}),
\end{aligned}
$$

# Hazard Function

Recall: $h(t) = -d(\log[S(t)])/dt$

$$h(x; Z) = h_0(\exp\{\mathbf{Z}\theta\}x)\exp(\mathbf{Z}\theta)$$

Example:

Suppose W follows an extreme value distribution with density function

$$f_W(w) = \exp\{w - e^w\}$$

Then:

$$S(x; Z) = S_0(\exp\{\mathbf{Z}\theta\}x) = \exp\{-[\exp\{\mathbf{Z}\theta\}x]^{1/\sigma}\}$$

$$h(x; Z) = h_0(\exp\{\mathbf{Z}\theta\}x) = (1/\sigma)[\exp\{\mathbf{Z}\theta\}x]^{1/\sigma-1}$$

where

$$S_0(t) = \exp(-t^{1/\sigma})$$

$$h_0(t) = (1/\sigma)\, t^{1/\sigma-1}$$

# Parametric Regression Model: exponential

```
proc lifereg data = unos_c;
     class yearsperf;
     model fu*death(0) = txtype age yearsperf /dist=exponential;
run;
```

```
                       The LIFEREG Procedure
                     Type III Analysis of Effects
                                        Wald
            Effect          DF     Chi-Square      Pr > ChiSq
            txtype           1       52.8909         <.0001
            age              1       29.7959         <.0001
            yearsperf        2        4.5849         0.1010
```

### Analysis of Maximum Likelihood Parameter Estimates

| Parameter | | DF | Estimate | Standard Error | 95% Confidence Limits | | Chi-Square | Pr > ChiSq |
|-----------|-----------|----|----------|----------|---------|---------|--------|--------|
| Intercept |           | 1  | 3.9121   | 0.2047   | 3.5109  | 4.3133  | 365.19 | <.0001 |
| txtype    |           | 1  | -0.7011  | 0.0964   | -0.8900 | -0.5121 | 52.89  | <.0001 |
| age       |           | 1  | 0.0470   | 0.0086   | 0.0301  | 0.0638  | 29.80  | <.0001 |
| yearsperf | 1990-1994 | 1  | 0.3885   | 0.1874   | 0.0212  | 0.7558  | 4.30   | 0.0382 |
| yearsperf | 1995-1999 | 1  | 0.3977   | 0.1929   | 0.0196  | 0.7757  | 4.25   | 0.0392 |
| yearsperf | 2000-2002 | 0  | 0.0000   | .        | .       | .       | .      | .      |
| Scale     |           | 0  | 1.0000   | 0.0000   | 1.0000  | 1.0000  |        |        |
| Weibull Shape |       | 0  | 1.0000   | 0.0000   | 1.0000  | 1.0000  |        |        |

```
                    Lagrange Multiplier Statistics
              Parameter       Chi-Square      Pr > ChiSq
              Scale            136.9235          <.0001
```

# Parametric Regression Model: exponential

```
. streg txtype age i.yearsperf, dist(exp)

------------------------------------------------------------------------------
         _t | Haz. Ratio   Std. Err.       z     P>|z|    [95% Conf. Interval]
------------+-----------------------------------------------------------------
     txtype |   2.045767    .1966766     7.45    0.000     1.694428    2.469956
        age |   .9536641    .0081702    -5.54    0.000     .9377845    .9698125
            |
  yearsperf |
       1995 |   .9991544    .0985305    -0.01    0.993     .8235544    1.212196
       2000 |    1.51563    .2801217     2.25    0.024     1.055049    2.177279
            |
      _cons |   .0135528    .0015876   -36.72    0.000     .0107725    .0170505
------------------------------------------------------------------------------

. streg txtype age ib(2000).yearsperf, dist(exp) nohr

------------------------------------------------------------------------------
         _t |      Coef.   Std. Err.       z     P>|z|    [95% Conf. Interval]
------------+-----------------------------------------------------------------
     txtype |   .7157727    .0961383     7.45    0.000     .5273451    .9042003
        age |  -.0474438    .0085672    -5.54    0.000    -.0642351   -.0306525
            |
  yearsperf |
       1990 |  -.4158315    .1848219    -2.25    0.024    -.7780758   -.0535872
       1995 |  -.4166775     .190232    -2.19    0.028    -.7895253   -.0438296
            |
      _cons |  -3.885333    .2022253   -19.21    0.000    -4.281687   -3.488979
------------------------------------------------------------------------------
```