Duration Analysis In Stata

Kevin Sweeney Assistant Director, Political Research Lab

Based On: An Introduction to Survival Analysis Using Stata

We Will Cover:

- 1. Overview Stata and "Shape" of Survival Data
- 2. ST-Setting and Describing Your Data
- 3. Nonparametric Analysis: Kaplan-Meier
- 4. Parametric Models (Exponential, Weibull...), and post-estimation
- 5. The Cox Proportional Hazards Model, and post-estimation

When You Open Stata...



ST Setting Your Data

he basic syntax is stset time_of_failure_or_censoring_variable, failure(one_if_failure_variable)

- So, if we had data that looked like this
- We'd type *stset failtime*

failtime	X
1	3
5	2
9	4
20	9
22	10

22	10	0	
20	9	1	
9	4	1	
5	2	1	
1	3	1	
lasttime	X	failed	

If we had data that looked like this, we'd type...

stset lasttime, failure(failed)

- _t0 & _t record time span
 _d records outcome
 _st records whether the
 observation is relevant.

ST Setting Your Data, Important Options

If	you have	more	than one	record	per subjec	ct you	ı must tell	Stata
wh	hat the id	varial	ole is	name	<u>lasttime</u>	X	failed	
tset lasttime, failure(failed) id(name)			Bob	1	3	1		
			Bob	5	2	1		
			Jim	9	4	1		
				Jim	20	9	1	
				Jim	22	10	0	
<u>name</u> Bob Bob Jim Jim Jim	<u>lasttime</u> 1 5 9 20 22	<u>x</u> 3 2 4 9 10	<u>event</u> 7 9 6 7 9	S	You can a kind of ev others are	also te vent is not , <i>failur</i>	ell Stata a s a failure, re(event==9	certain , wherea 9) <i>id(name</i>
			1					

ST Setting Your Data, One More Option

Finally (well, not really), you can tell Stata when your observations begin, if you don't Stata will do it for you...

<u>name</u>	<u>begin</u>	<u>lasttime</u>	X	event
Bob	0	1	3	7
Bob	3	5	2	9
Jim	0	9	4	6
Jim	17	20	9	7
Jim	21	22	10	9

stset lasttime, failure(event==9) id(name) time0(begin)

Now, on to a real live example...

Canned Hip Fracture Data

use http://www.stata-press.com/data/cgg/hip2

Type: *describe* to see what you have

С	Contains data	from http	p://www.s	tata-press.com	n/data/cgg/hip2.dta	
	obs:	106			hip fracture study	
	vars:	12			30 Jan 2002 17:58	
	size:	2,332 (96.7% of 1	memory free)		
-	variable name	storage type	display format	value label	variable label	 \
- i t f f a c	d ime0 ime1 fracture protect ge alcium	byte byte byte byte byte byte float	<pre>%4.0g %5.0g %5.0g %8.0g %8.0g %4.0g %8.0g %8.0g</pre>		patient id begin of span end of span fracture event wears device age at enrollment blood calcium level	

Canned Hip Fracture Data



ST Setting the Data

Type: *stset time1, id(id) time0(time0) failure(fracture)* And this is what you get...

> id: id failure event: fracture ~= 0 & fracture ~= . obs. time interval: (time0, time1] exit on or before: failure _____ 106 total obs. 0 exclusions 106 obs. remaining, representing subjects 48 failures in single failure-per-subject data 31 714 total analysis time at risk, at risk from t = 0earliest observed entry t = 0

> > last observed exit t = 39

This may not make a lot of sense, Stata has a more descriptive command...

stdes

Type: *stdes* and we see that...



Nonparametric Analysis: Kaplan-Meier

$$\hat{S}(t) = \prod_{j|t_j \le t} \left(\frac{n_j - d_j}{n_j}\right)$$

-the probability of survival past time *t*, or the probability of failing after time *t*.

Where n_j is the number of individuals at risk at t_j and d_j is the number of failures at t_j .

Type sts graph

And you get the simple Kaplan-Meier graph



STS List

🗰 Intercooled Stata 7.0			_ 🗆 ×
<u>File E</u> dit <u>P</u> refs <u>W</u> indow <u>H</u> elp			
	🛛 🚿 🔲 🙆 💿	8	
🚍 Stata Results			×
failure_d: fracture analysis time_t: time1 id: id			
Beg. Net Time Total Fail Lost	Survivor Std. Function Error	[95% Conf. Int.]	
1 48 2 0 2 46 1 0 3 45 1 0 5 42 2 3 6 37 2 1 0 5 42 2 3 6 37 2 1 0 8 33 3 1 9 29 0 1 10 28 1 1 12 23 2 0 13 21 1 0 15 20 1 -2 16 21 1 0 15 20 1 -2 16 21 1 0 19 18 0 2 15 20 1 -2 16 21 1 0 19 18 0 1 22 15 2 0 16 0 1 22 15 2 0 24 11 1 17 20 1 0 22 15 2 0 18 0 1 19 18 0 1 22 16 0 1 23 4 1 0 19 18 0 1 23 18 0 1 24 10 1 25 20 1 0 10 2 8 1 0 10 2 8 1 0 10 2 8 1 0 11 1 0 15 20 1 0 15 20 1 0 15 20 0 16 0 0 1 22 1 1 0 23 18 0 0 24 10 1 0 24 10 1 0 25 20 1 0 10 1 0 27 10 1 0 28 8 8 0 0 20 1 0 29 10 1 0 20 1 0 20 1 0 20 1 0 20 1 0 20 1 0 21 0 0 22 0 0 10 1 0 23 18 0 0 24 10 0 24 10 0 25 20 0 10 1 0 27 10 0 28 8 0 0 20 1 0 29 10 1 0 20 0 1 31 0 20 1 0 20 10 0	$\begin{array}{ccccccc} 0.9583 & 0.0288 \\ 0.9375 & 0.0349 \\ 0.9167 & 0.0399 \\ 0.8750 & 0.0477 \\ 0.8333 & 0.0538 \\ 0.7883 & 0.0538 \\ 0.7651 & 0.0683 \\ 0.6955 & 0.0683 \\ 0.6955 & 0.0683 \\ 0.6955 & 0.0683 \\ 0.6955 & 0.0737 \\ 0.5653 & 0.0765 \\ 0.5384 & 0.0774 \\ 0.5114 & 0.0781 \\ 0.4627 & 0.0779 \\ 0.4627 & 0.0778 \\ 0.3085 & 0.0766 \\ 0.2776 & 0.0749 \\ 0.2429 & 0.0731 \\ 0.2429 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1822 & 0.0760 \\ 0.1824 & 0.0760 \\$	0.8435 0.9894 0.8186 0.9794 0.7930 0.9679 0.7427 0.9418 0.6943 0.9129 0.6419 0.8802 0.6155 0.8627 0.5397 0.8076 0.5397 0.8076 0.5397 0.8076 0.5130 0.7874 0.4585 0.7447 0.4033 0.6988 0.3767 0.6511 0.3285 0.6283 0.3066 0.6052 0.3066 0.6052 0.3487 0.0638 0.3487 0.0638 0.3487	

Nonparametric Analysis: Kaplan-Meier

We can make the graph a little more complicated by comparing those in the treatment group with those in the control group.

Use the *by* command to plot multiple survival curves...

Type sts graph, by(protect)



Those without the protection fail more quickly than those with it.

Simple Nonparametric Tests



Parametric Models - Exponential

Type: <u>streg</u> <u>age protect</u>, <u>dist(exp)</u> <u>nohr</u>

- The command for all Parametric Models
- The covariates in this model.
- Specify the parameterization of the baseline hazard.
- Tells Stata you want coefficients and not hazard ratios => must exponentiate.

$$h(t) = \exp(\mathbf{b}_0 + \mathbf{b}'_k X_i)$$

 $(t) = \exp(-7.89) = .00037$ $(t | x_j) = .00037 \exp(-1.69 \text{ protect} + .08 \text{ age}_j)$

intercooled	Stata 7.0						۱×
<u>File E</u> dit <u>P</u> refs	<u>W</u> indow <u>H</u>	elp					
F	S: 👁		Š		0 😣		
🔚 Stata Resu	lts						
. streg age pr	otect, dist(exp) nohr					
failu analysis ti	me_d: frac me_t: time id: id	ture 1					
Iteration 0: Iteration 1: Iteration 2: Iteration 3: Iteration 4: Iteration 5:	log likelih log likelih log likelih log likelih log likelih log likelih	ood = -60.06 ood = -54.03 ood = -47.55 ood = -47.53 ood = -47.53 ood = -47.53	7085 4598 3588 4671 4656 4656				
Exponential re	gression	log relative	-hazard (Form			
No. of subject No. of failure Time at risk	S = S =	48 31 714		Numb	er of obs	- 1	06
Log likelihood	47.53	4656		LR c Prob	hi2(<mark>2</mark>) ⇒_chi2	= 25. = 0.00	.06 000
_t	Coef.	Std. Err.	z	P>1z1	E95% Cor	nf. Interva	at 3
protect _cons	.0809663 -1.688958 -7.892737	.0342787 .3703357 2.458841	2.36 -4.56 -3.21	0.018 0.000 0.001	.0137813 -2.414803 -12.71199	3 .14819 396311 8 -3.0734	514 137 198
CADOTA							-

Parametric Models - Exponential

Or, you could estimate the model and get hazard ratios...

Type: *streg age protect, dist(exp)*

Remember the coefficient on age was .0809663, $e^{.0809663}=1.084334$

Hazard ratios have the virtue of being relatively easy to interpret.

🚅 Intercooled Stata 7.0					_ 🗆 🗵
<u>File E</u> dit <u>P</u> refs <u>W</u> indow	<u>H</u> elp				
	> 🖂 🔜 🖉	3		o 😣	
🚍 Stata Results					2
. streg age protect, dis	t(exp)				
failure_d: fr analysis time_t: ti id: id	acture mel				
Iteration 0: log likel Iteration 1: log likel Iteration 2: log likel Iteration 3: log likel Iteration 4: log likel Iteration 5: log likel	ihood = -60.06 ihood = -54.03 ihood = -47.55 ihood = -47.53 ihood = -47.53 ihood = -47.53	7085 4598 3588 4671 4656 4656			
Exponential regression -	- log relative	-hazard	form		
No. of subjects = No. of failures =	48 31		Number	r of obs =	106
Time at risk = Log likelihood = -47.	714 534656		LR ch Prob	i2(2) = > chi2 =	25.06 0.0000
_t Haz. Rati	o Std. Err.	z	P>1z1	E95% Conf.	Interval]
age 1.08433 protect .184711	4 .0371696 8 .0684054	2.36 -4.56	0.018 0.000	1.013877 .0893849	1.159688 .3817025
C:\DATA					

Parametric Models - Weibull

Type: streg age protect, dist(weib) nohr

Here, Stata estimates the shape of the hazard function with p.

p>1 indicates the hazard is monotonically increasing.

p < 1 indicates the hazard is monotonically decreasing.

Gompertz, lognormal, loglogistic, gamma

	📕 Intercooled	Stata 7.0					- 🗆 ×
	<u>File E</u> dit <u>P</u> ref:	s <u>W</u> indow <u>H</u> e	lp				
	-	51 T	.	3		0 8	
1	🚍 Stata Res	ults					
	. streg age p	rotect, dist()	weib) nohr				
	fail analysis t	ure_d: frac ime_t: time id: id	ture 1				
	Fitting const	ant-only mode	L:				
	Iteration O: Iteration 1: Iteration 2: Iteration 3:	log likelih log likelih log likelih log likelih	ood = -60.06 ood = -59.3 ood = -59.29 ood = -59.29	7085 0148 8481 8481			
	Fitting full	model:					
	Iteration 0: Iteration 1: Iteration 2: Iteration 3: Iteration 4: Iteration 5: Weigull regre	log likelih log likelih log likelih log likelih log likelih log likelih ssion log :	ood = -59.29 ood = -54.88 ood = -42.12 ood = -41.99 ood = -41.99 ood = -41.99 relative-haz	8481 7563 3875 3012 2704 2704 2704 ard form			
	No. of subjec	ts =	48		Numb	er of obs	= 106
	Ho. of failur Time at risk Log likelihoo	es = = d = -41.99	714 2704		LR c Prob	hi2(<mark>2</mark>) > chi2	= 34.61 = 0.0000
	_	Coef.	Std. Err.	z	P>(z)	E95% Conf	F. Intervall
	age protect _cons	.1108134 -2.207628 11.67104	.0378734 .4076113 2.90919	2.93 -5.42 -4.01	0.003 0.000 0.000	.036583 -3.006532 -17.37295	.1850439 -1.408725 -5.969135
	∕ľn_p	.5188694	.1376486	3.77	0.000	.2490831	.7886556
	1∕p	1.680127	.2312671 .0819275			1.282849 .4544553	2.200436 .7795152
0	CADATA						

Parametric Models, postestimation

median time	predicted median survival time; the default					
median lntime	predicted median ln(survival time)					
mean time	predicted mean survival time					
mean lntime	predicted mean ln(survival time)					
hazard	predicted hazard					
hr	predicted hazard ratio					
xb	linear prediction					
stdp	standard error of the linear prediction					
surv	predicted S(depvar) or S(depvar t0)					
csnell	(partial) Cox-Snell residuals					
mgale	(partial) martingale-like residuals					
deviance	deviance residuals					
csurv	predicted S(depvar earliest t0 for subject)					
ccsnell	cumulative Cox-Snell residuals					
cmgale	cumulative martingale-like residuals					

A Post-estimation Example

After our most recent regression (the Weibull) we could Type: <u>predict cc</u>, <u>ccsnell</u> <u>graph cc _t</u>, <u>s([id])</u>



We will consider some more post-estimation commands with the Cox Model

The Semiparametric Cox Proportional Hazards Model $h(t | x) = h_0(t) \exp(\mathbf{b}'_k x_i)$

Type: stcox age protect, nohr



Extensions to the Cox Model

Type: *stcox age protect, nohr <u>basechazard(HO)</u>*



🗱 Intercooled S	Stata 7.0					_ 🗆 ×
<u>File E</u> dit <u>P</u> refs	<u>W</u> indow <u>H</u> e	lp				
🖻 🖬 🎒	S: 👁		3		00	
🚍 Stata Resul	ts					
. stcox age pro	otect, nohr b	asechazard()	HO)			
failur analysis tim	ne _d: fract ne _t: time1 id: id	ure				
Iteration 0: Iteration 1: Iteration 2: Iteration 3: Refining estima Iteration 0:	log likeliho log likeliho log likeliho log likeliho ates: log likeliho	d = -98.57 d = -82.739 d = -82.47 d = -82.47 d = -82.47 d = -82.47	1254 5029 1037 0259 0259			
Cox regression	Breslow m	method for t	les:			
No. of subjects No. of failures	5 = 5 =	48 _31		Numbe	er of obs	= 106
Time at risk Log likelihood	= = -82.470	714 0259		LR cl Prob	ni2(2) ≻ chi2	= 32.20 = 0.0000
_t _d	Coef.	Std. Err.	z	P>(z)	E95% Conf	. Interval]
age protect	.1052352 -2.256836	.0378119 .4538632	2.78 -4.97	0.005	.0311253 -3.146392	.1793451 -1.367281
C:\DATA						

Type: graph HO _t, c(J) sort

Testing the Proportional Hazards Assumption: Schoenfeld Residuals

Type: *stcox age protect*, <u>*schoenfeld(sch*)*</u>) <u>*scaledsch(sca*)*</u>

After estimation type: *stphtest, detail*

🐖 Intercooled Stata	7.0				_ 🗆 ×
<u>File Edit Prefs Wind</u>	ow <u>H</u> elp				
B B S		. 🧕 🗆		0 😣	
🚍 Stata Results					×
. stphtest, detail	190000000000000000				
Time: Time	tional hazards as	ssumption			
True, True	rho	chi2	df	Prob>chi2	
age protect	-0.11519 0.00889	0.43 0.00	1	0.5140 0.9627	
global test		0.44	2	0.8043	
C:\DATA					

We find no evidence that the model violates the PH assumption.

Testing the Proportional Hazards Assumption: Stata's Plots: *stphplot*

stphplot estimates $-\ln[-\ln{S(t)}]$ vs. $\ln(t)$ for each level of the specified covariate.



Type:stphplot, by(protect) adjust(age



Parallel lines means model has not violated PH assumption.