Class 7 Outline

- 1. Review of Cox proportional hazards regression model using a detailed example
- Risk Factors for Infant Mortality in Sarlahi
 District, Nepal: A Survival Analysis of Data
 from the Nepal Nutrition Intervention Project
 (NNIPS-II)

1

0. Learning Objectives

- Explore a large survival data set, identifying data limitations, and data patterns useful in guiding decisions about a statistical analysis
- Use Kaplan-Meier survival curve estimates and crude survival rates to study dependence of survival on discrete characteristics of the mother and pregnancy
- Create predictor variables as candidates for survival regression models
- Choose, fit, and check Cox proportional hazards models using time independent predictor variables

1. Scientific Background

 NNIPS-II (Nepal Nutrition Intervention Project Sarlahi- II) was a randomized study of 15,987 infants born alive to 43,559 women who received Vitamin A, Beta-carotene, or a placebo prior to and during pregnancy

3

1.1 Objectives of Analysis

- Investigate effect of treatment on infant mortality
- Identify maternal and pregnancy characteristics that predict infant mortality
- Estimate shape of infant mortality rate curve as a function of gestational age and parity
- Estimate multivariable Cox model to predict survival as a function of treatment and other mother/child characteristics such as gestational age and parity

1.2 Data Exploration and "Cleaning"

- We are using a representative subset of 10,295 live births from the original data set
- Women/infant pairs followed for up to 6 months (180 days) post delivery

5

1.3a Key Variables

- Survival time (stime) in days
- Censoring indicator (cens): 1 if death, 0 if censored
- Mother's parity (parity): number of previous live births

1.3b Key Variables

- Night blindness (nblind): 1 if mother night blind during pregnancy; 0 if not
- Gestational Age (gestage): number of weeks from last completed menstrual cycle to birth
- Treatment (treat): 1 if Beta-carotene, 2 if placebo, 3 if Vitamin A
- Gender (male): 1 if male, 0 female

7

1.4a Codebook, Key Variables

1.4b Codebook, Key Variables

```
. codebook parity nblind
```

1.4c Codebook, Key Variables

. codebook gestage treat

1.4d Codebook, Key Variables

2.0 Create "Survival" Data Set

```
. stset stime, failure(cens=1)

failure event: cens == 1
obs. time interval: (0, stime)
exit on or before: failure

10295 total obs.

518 obs. end on or before enter()

9777 obs. remaining, representing
559 failures in single record/single failure data
1704029 total analysis time at risk, at risk from t = 0
earliest observed entry t = 0
last observed exit t = 918
```

2.1a Why have Observations been Dropped?

 518 observations have just been dropped by Stata because survival time (*stime*) recorded as 0 days

. tab cens if stime == 0

Cum.	Percent	Freq.	cens
74.13 100.00	74.13 25.87	384 134	0
	100.00	518	Total

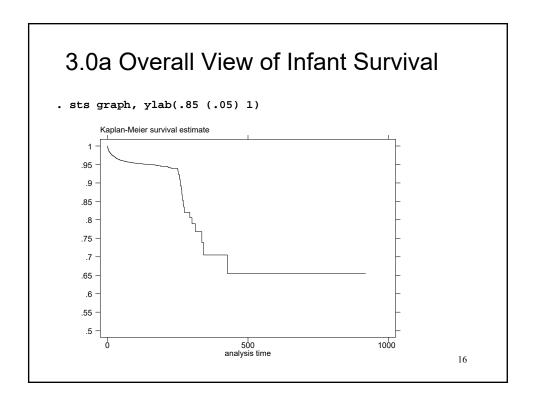
· Infants who died on day 1

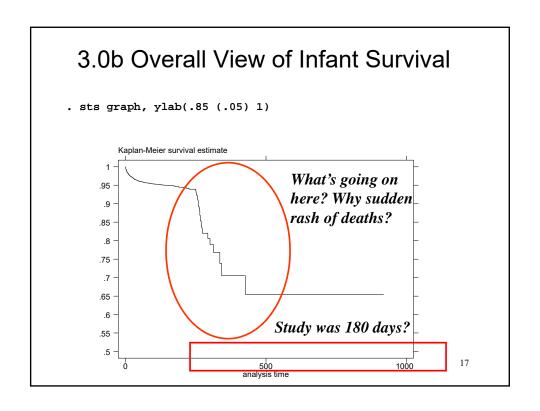
13

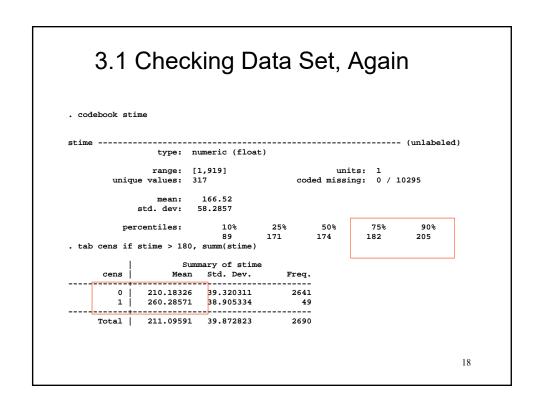
2.1b Why have Observations been Dropped?

- These 518 observations give information about survival experience!
- Remedy: add 1 to all survival times so that Stata will use these observations in its computations

2.2 Remedy for Dropped Observations







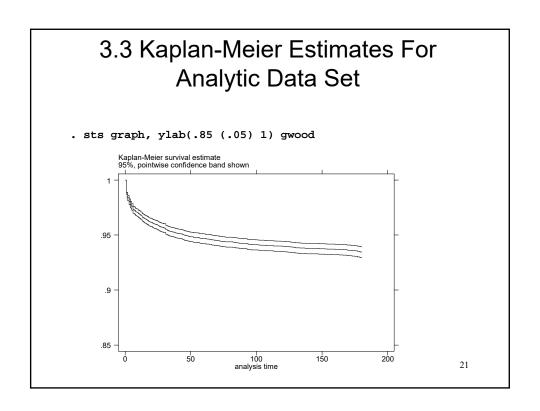
3.2a Follow-up With Pls

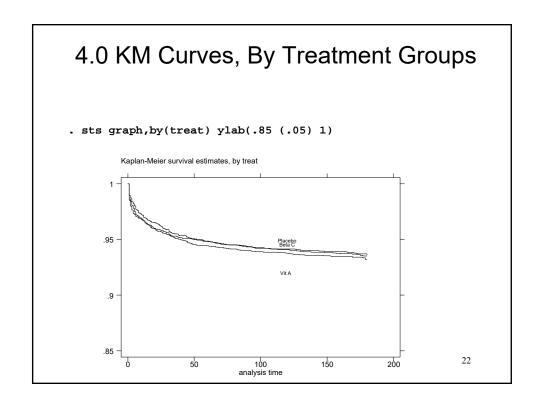
- According to the PIs, some women were followed beyond 180 days in "passive surveillance"
- Formal study ("active surveillance") ended at 180 days even though additional data was collected
- Selection bias in using time to event information for this subset of women

19

3.2b Follow-up With Pls

"Remedy" – censor all observations at 180 days





4.1 Incidence Rates, By Treatment Groups

. stsum, by(treat)

failure _d: cens == 1
analysis time _t: stime

treat	 time at risk	incidence rate	no. of subjects	Sur 25%	vival time 50%	 75%
Beta C	+ 516692	.0003929	3265	•		
Placebo	532438	.000385	3387	•	•	
Vit A	578595	.0004079	3643	•	•	•
total	+ 1627725	.0003956	10295	···		

23

4.2 Hazard Ratios for Treatment Groups

5.0 Gestational Age

First, create some categories for *gestage*

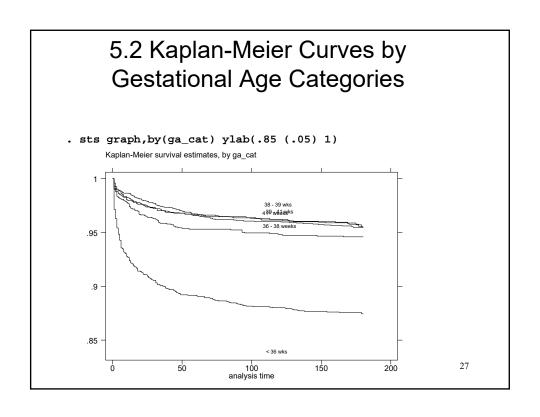
```
.gen ga_cat=1 if gestage < 36 & gestage ~=.
.replace ga_cat=2 if gestage >=36 & gestgage < 38
.replace ga_cat=3 if gestage >=38 & gestgage < 39
.replace ga_cat=4 if gestage >=39 & gestgage < 41
.replace ga_cat=5 if gestage >=41
.label define gestcat 1 "< 36 wks" 2 "36-38 weeks" 3 "38-39 weeks" 4 "39-41 weeks" 5 "41+ weeks"</pre>
```

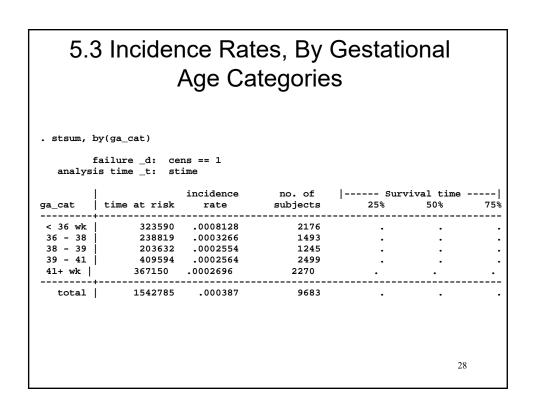
24

5.1 Distribution of Gestational Age Categories

. tab ga_cat

ga_cat	Freq.	Percent	Cum.
< 36 wks 36 - 38 weeks 38 - 39 wks 39 - 41 wks 41+ weeks	2176 1493 1245 2499 2270	22.47 15.42 12.86 25.81 23.44	22.47 37.89 50.75 76.56 100.00
+ Total	9683	100.00	





5.4 Incidence Rates by Gestational Age (With Missing Values)

. stsum, by(ga_cat_miss)

failure _d: cens == 1
analysis time _t: stime

		incidence	no. of	Surv	vival time	
ga_cat~s	time at risk	rate	subjects	25%	50%	75%
< 36 wk	+ 323590	.0008128	2176			
36 - 38	238819	.0003266	1493	•	:	:
38 - 39	203632	.0002554	1245	•	•	
39 - 41	409594	.0002564	2499	•	•	
41+ wk	367150	.0002696	2270	•	•	
Missing	84940	.0005533	612		•	•
total	1627725	.0003956	10295	•	•	•

Note high mortality among those missing gestational age!

29

6.0 Mother's Parity Categories

First create categories for parity

```
.gen par_cat=0 if parity ==0
```

- .replace par_cat=1 if parity==1
- .replace par_cat=2 if parity >=2 & parity < =4</pre>
- .replace par_cat=3 if parity >=5 & parity < =8</pre>
- .replace par_cat=4 if parity > 8 & parity~=.
- .label define par 0 "No prev child" 1 "1 prev child" 2 "2-4 prev child" 3 "5-8 prev child" 4 "8+ prev child"
- .label values par_cat par

6.1 Distribution of Parity Categories

. tab par_cat

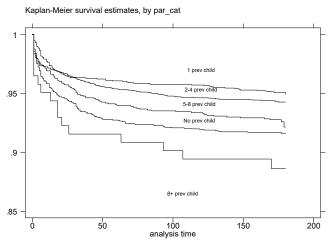
par_cat	Freq.	Percent	Cum.
No prev child 1 prev child 2-4 prev child 5-8 prev child 8+ prev child	2254 2018 4262 1363 142	22.45 20.10 42.45 13.58 1.41	22.45 42.55 85.01 98.59 100.00
Total	10039	100.00	

3

32

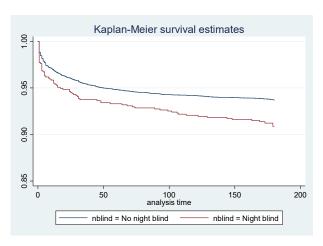
6.2 Kaplan-Meier Curves by Parity Categories

. sts graph,by(par_cat) ylab(.85 (.05) 1)



7.0 Kaplan-Meier Curves, By Mother's Night Blindness

. sts graph,by(nblind) ylab(.85 (.05) 1)



33

7.1 Incidence Rates, By Mother's Night Blindness

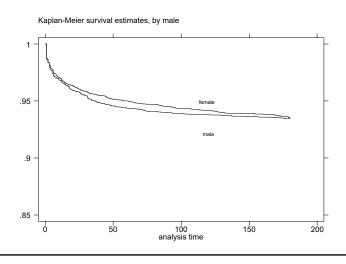
. stsum, by(nblind)

failure _d: cens == 1
analysis time _t: stime

	1	incidence	no. of	Surv	vival time	
nblind	time at risk	rate	subjects	25%	50%	75%
No night	1482249	.0003805	9372	•	•	
Night bl	145476	.0005499	923			
total	1627725	.0003956	10295			

8.0 Kaplan-Meier Curves, By Gender (excluding missing)

. sts graph, by(male) ylab(.85 (.05) 1)



8.1 Incidence Rates, By Gender

. stsum, by(male)

failure _d: cens
analysis time _t: stime

male	time at risk	incidence rate	no. of subjects	Surv 25%	vival time 50%	 75%
female male	840324 872832	.0003689	4966 5195	· ·	· ·	:
total	1713156	.0003742	10161	•		

9.0a Cox Proportional Hazards Model

- h(t,X) = hazard for people at risk with predictor values X = (X₁,X₂,X_p)
- · We can write

$$h(t,X) = h_0(t) \cdot e^{\beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p}$$

• Similarly, $log[h(t,\mathbf{x})] = log[h_o(t)] + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$

37

9.0b Cox Proportional Hazards Model

• Hazard Ratio (relative hazard) interpretation of the $\beta\ensuremath{\text{'s}}$

$$e^{\beta_1} = \frac{h(t, X_1 + 1, X_2, ... X_p)}{h(t, X_1, X_2, ... X_p)}$$

= relative hazard by one unit increase in X_1 , but have same values for X_2 , X_p (at any fixed time t)

9.1 Interpretation of β 's

• Why do we have a hazard ratio interpretation?

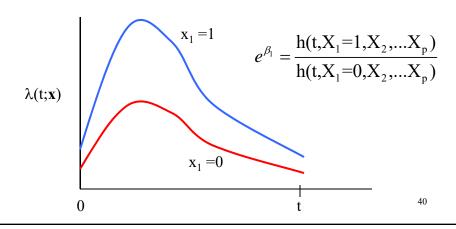
$$\begin{split} &\log[h(t, X_1 + 1, X_2, ... X_p)] = \log[h_o(t)] + \beta_1(X_1 - 1) + \beta_2 X_2 + ... \beta_p X_p \\ &- \log[h(t, X_1, X_2, ... X_p)] = \log[h_o(t)] + \beta_1 X_1 + \beta_2 X_2 + ... \beta_p X_p \end{split}$$

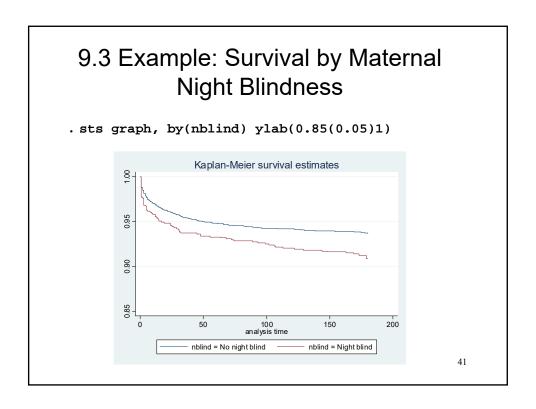
$$\beta_1 = \log \left[\frac{h(t, X_1 + 1, X_2, ... X_p)}{h(t, X_1, X_2, ... X_p)} \right]$$

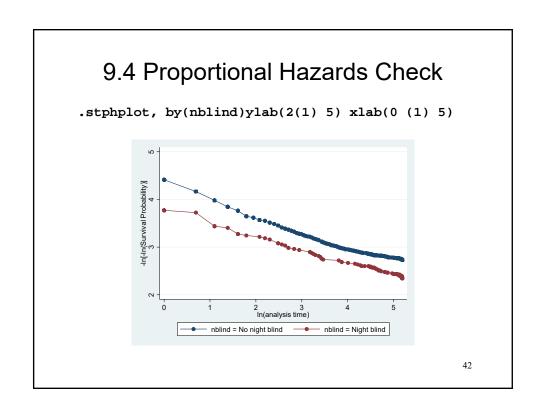
39

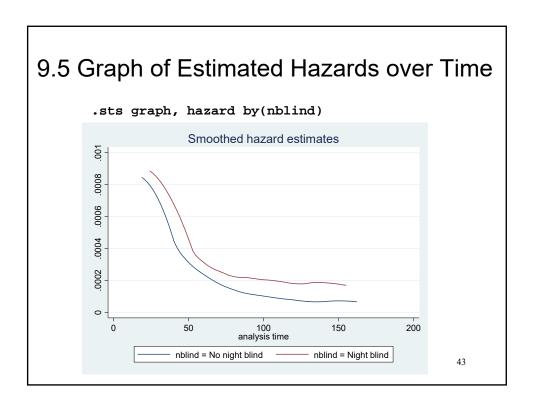
9.2 Cox Proportional Hazards Model Assumption

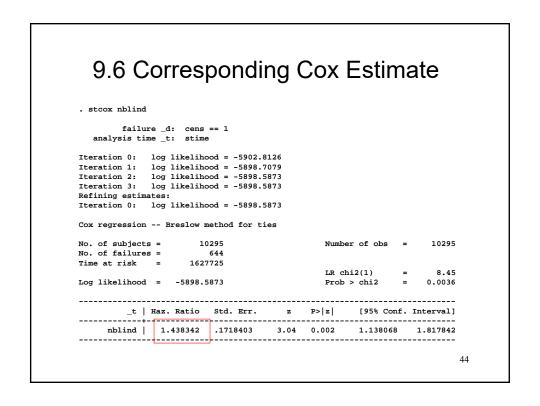
· Assumes proportional hazards over time











9.7 Interpretation

- Children of night-blind mothers are estimated to be at higher risk of death than children of non-night-blind mothers, averaged over the study duration
- The estimated relative hazard of death for children of night-blind versus non-night-blind mothers is 1.44 (95% CI 1.14 – 1.82) at any time in the 6-month follow-up period

45

10. Analysis Plan

- Using Cox Regression
 - First, look at gestational age and parity individually to choose a model form
 - Second, build a multivariable model for treatment including gestational age and parity, as well as other predictors (gender, mother's night blindness)

10.1 Gestational Age, Categorical

. tab ga_cat

ga_cat	Freq.	Percent	Cum.
<pre></pre>	2176 1493 1245 2499 2270	22.47 15.42 12.86 25.81 23.44	22.47 37.89 50.75 76.56 100.00
Total	9683	100.00	

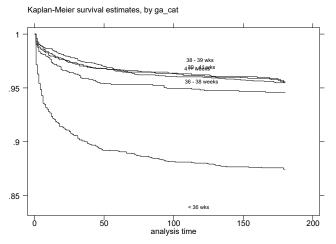
Recall, total n = 10,295: 612 observations missing gestational age

4

48

10.2 Kaplan-Meier Curves by Gestational Age Categories

. sts graph, by(ga_cat) ylab(.85 (.05) 1)



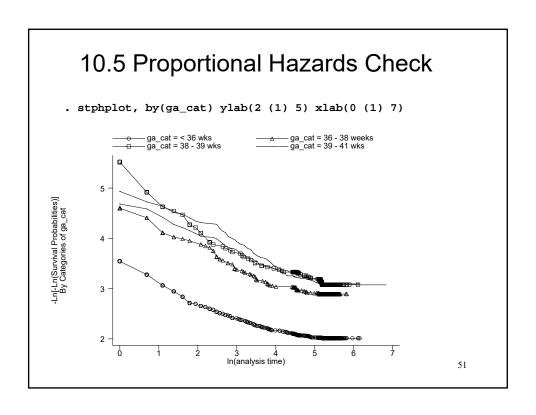
10.3 Cox Results with Hazard Ratios (HRs)

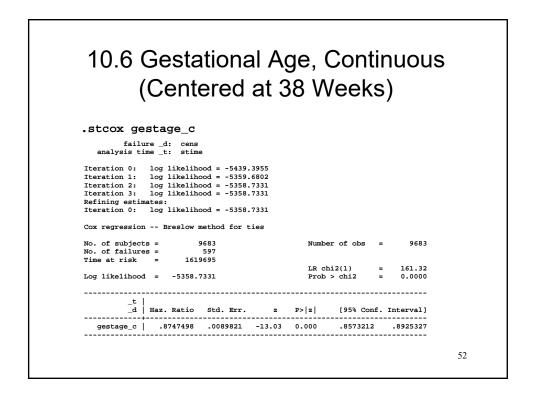
40

10.4 Cox Results with log HRs

```
.stcox i.ga_cat, nohr
```

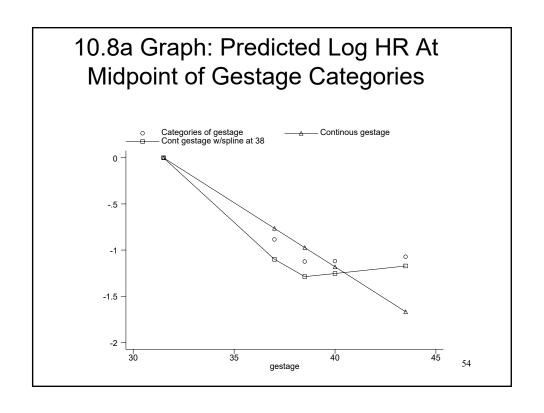
```
_Iga_cat_1-5
                                                    (naturally coded; _Iga_cat_1 omitted)
              failure _d: cens == 1
analysis time _t: stime
           Iteration 0: log likelihood = -5439.3955
Iteration 1: log likelihood = -5368.3146
Iteration 2: log likelihood = -5362.5375
Iteration 3: log likelihood = -5362.5358
Refining estimates:
            Iteration 0: log likelihood = -5362.5354
            Cox regression -- Breslow method for ties
                                      9683
           No. of subjects = No. of failures =
                                                                   Number of obs =
                                                                                             9683
                                  1542785
            Time at risk =
                                                                   LR chi2(4)
                      _t |
_d |
           __d | Coef. Std. Err. z P>
                                                                           [95% Conf. Interval]
```





10.7 Add a Spline Term for Gestational Age at 38 Weeks

```
.stcox gestage_c ga_sp38
    failure _d: cens == 1
analysis time _t: stime
                     log likelihood = -5439.3955
Iteration 0:    log likelihood = -5439.3955
Iteration 1:    log likelihood = -5375.191
Iteration 2:    log likelihood = -5338.943
Iteration 3:    log likelihood = -5338.7324
Iteration 4:    log likelihood = -5338.7324
Iteration 0: log likelihood = -5338.7324
Cox regression -- Breslow method for ties
No. of subjects = No. of failures =
                                     9683
                                                                        Number of obs
                                                                                                         9683
Time at risk =
                                                                        LR chi2(2)
                                                                                                       201.33
Log likelihood = -5338.7324
                                                                                                       0.0000
                                                                        Prob > chi2
             _t |
_d | Haz. Ratio Std. Err.
                                                                                  [95% Conf. Interval]
    gestage_c |
                        .8185724 .0115865 -14.14 0.000
1.258339 .0440248 6.57 0.000
                                                                                   .7961754
                                                                                                    .8415995
                     1.258339
                                                                               1.174944
       ga_sp38
```



10.8b Notes on Graph of Predicted log(HR)

- Estimated log HR of death for children in each gestational age group as compared to children in lowest gestational age group
- Both the spline model and the model with gestational in categories yield similar estimates, and same trend

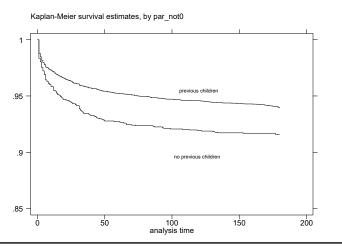
55

11.0 Parity, As Binary Categorical

```
. codebook parity
type: numeric (float)
               range: [0,15]
values: 15
                                            units: 1
missing .: 256/10,295
        unique values: 15
             mean: 2.32563
std. dev: 2.14178
          percentiles: 10%
. gen par_not0=1 if parity >0 & parity < 16
(2,510 missing values generated)
. replace par_not0=0 if parity==0
(2,254 real changes made)
. tab par_not0
  par_not0 | Freq. Percent
  0 | 2,254 22.45 22.45
1 | 7,785 77.55 100.00
     Total | 10,039 100.00
                                                                       56
```

11.1 Kaplan-Meier Curves by Parity (Binary)

. sts graph, by(par_not0) ylab(.85 (.05) 1)



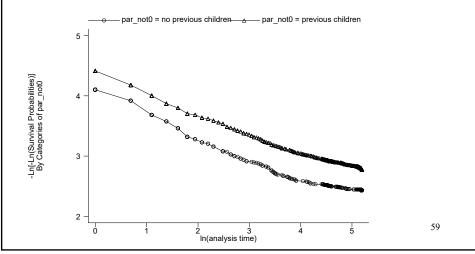
11.2 Cox Results with Parity, As Binary Categorical

.stcox par_not0

```
failure _d: cens
analysis time _t: stime
Iteration 0: log likelihood = -5781.046
Iteration 1: log likelihood = -5773.349
Iteration 2: log likelihood = -5773.268
Iteration 3: log likelihood = -5773.268
Refining estimates:
Iteration 0: log likelihood = -5773.268
Cox regression -- Breslow method for ties
No. of subjects =
                                10039
                                                                  Number of obs
                                                                                              10039
No. of failures =
                             1683476
Time at risk =
                                                                                              15.56
                                                                  LR chi2(1)
Log likelihood =
                         -5773.268
                                                                  Prob > chi2
                                                                                             0.0001
            _t |
_d | Haz. Ratio Std. Err.
                                                     z \qquad P > |z|
                                                                        [95% Conf. Interval]
     par_not0 | .6995863 .0616632 -4.05 0.000
                                                                        .588592 .8315114
```



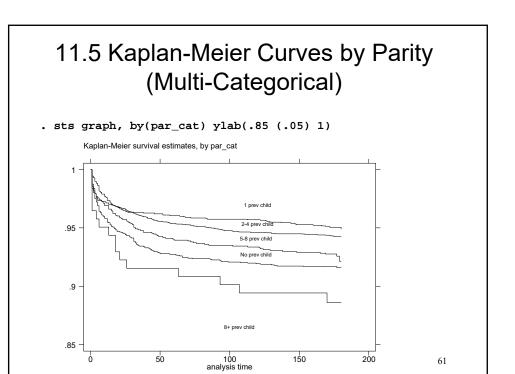
. stphplot, by(par_not0) ylab(2 (1) 5) xlab(0 (1) 5)

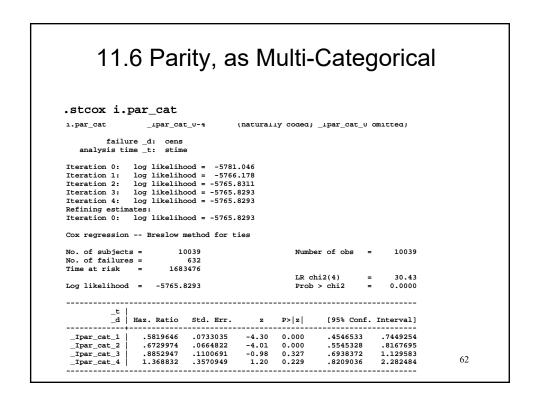


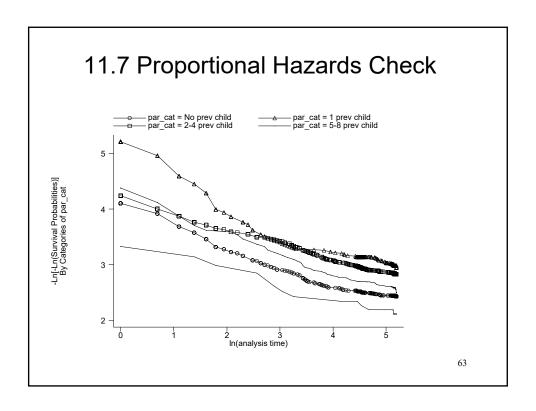
11.4 Parity, as Multi-Categorical

. tab par_cat

par_cat	Freq.	Percent	Cum.
No prev child	2254	22.45	22.45
1 prev child	2018	20.10	42.55
2-4 prev child	4262	42.45	85.01
5-8 prev child	1363	13.58	98.59
8+ prev child	142	1.41	100.00
Total	10039	100.00	







11.8 How to Model the Mortality and Parity Relationship?

- Both the binary and categorical modeling of parity showed strong parity/mortality association
- Categorical approach allowed for better modeling of relationship
- We could have also used a linear predictor approach (and also examined spline terms)

12.0 A Multivariable Model for Describing Child Mortality

- Two purposes
 - Estimating adjusted relationships between hazard of mortality and each predictor
 - Building a model to estimate hazard (and subsequently survival) of children with certain child and/or maternal characteristics

65

12.1a Cox Regression Results (HRs)

.stcox i.ga_cat i.par_cat i.male i.nblind i.treat if male~=9

```
_Iga_cat_1-5 (naturally coded; _Iga_cat_1 omitted)
_Ipar_cat_0-4 (naturally coded; _Ipar_cat_0 omitted)
_Imale_0-9 (naturally coded; _Imale_0 omitted)
_Inblind_0-1 (naturally coded; _Inblind_0 omitted)
_Itreat_1-3 (naturally coded; _Itreat_1 omitted)
i.ga cat
i.par cat
i.nblind
i.treat
                failure _d: cens == 1
      analysis time _t: stime
note: _Imale_9 dropped due to collinearity
note: _Imale_9 dropped due to collinearity
Iteration 0: log likelihood = -5312.0553
Iteration 1: log likelihood = -5222.0508
Iteration 2: log likelihood = -5215.5804
Iteration 4: log likelihood = -5215.576
Iteration 4: log likelihood = -5215.576
Refining estimates:
Iteration 0: log likelihood = -5215.576
Cox regression -- Breslow method for ties
No. of subjects =
                                                                                                Number of obs =
                                                                                                                                             9443
No. of failures =
                                       1523838
Time at risk =
                                                                                                LR chi2(12)
                                                                                                                                         192.96
Log likelihood = -5215.576
                                                                                                                                         0.0000
```

12.1b Cox Regression Results (HRs)

_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf.	<pre>Interval]</pre>
_Iga_cat_2	.4118807	.0537904	-6.79	0.000	.3188651	.5320296
_Iga_cat_3	.3230919	.0495576	-7.37	0.000	.2392017	.4364033
_Iga_cat_4	.3212073	.0377706	-9.66	0.000	.2550898	.4044619
_Iga_cat_5	.3469316	.0413817	-8.88	0.000	.2746082	.438303
_Ipar_cat_1	.5396267	.0716611	-4.65	0.000	.4159641	.7000531
_Ipar_cat_2	.640755	.0654042	-4.36	0.000	.524574	.7826674
_Ipar_cat_3	.7880633	.1013753	-1.85	0.064	.6124403	1.014048
_Ipar_cat_4	1.179086	.3282661	0.59	0.554	.683227	2.034821
_Imale_1	1.008702	.0836094	0.10	0.917	.8574488	1.186635
_Inblind_1	1.407019	.1767372	2.72	0.007	1.099966	1.799783
_Itreat_2	.9556154	.0988269	-0.44	0.661	.7802871	1.170339
_Itreat_3	.9609895	.0968532	-0.39	0.693	.7887337	1.170865

67

12.2 Summary of Multivariable Model Results

- Notice that hazard ratio estimates for both gestational age and parity categories are similar to the estimates obtained in the simple models (Why is this?)
- Each hazard ratio estimate is interpretable as an adjusted hazard ratio – interpretation depends on coding of predictor

12.3a Summary of Multivariable Model Results

- · Example: Night blindness
 - Coefficient estimates the adjusted log HR for children with night-blind mothers compared to children with non-night-blind mothers

69

70

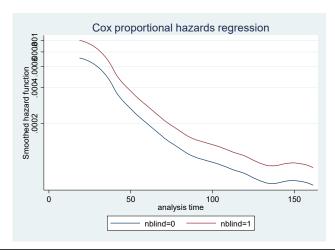
12.4 Adjusted Kaplan-Meier Curves by Maternal Night Blindness After .stcox for the multivariable model: .stcurve, survival at1(nblind=0) at2(nblind=1) ylab(0.85 (0.05) 1) Cox proportional hazards regression Cox proportional hazards regression

nblind=0

nblind=1

12.5 Adjusted Hazard Curves by Maternal Night Blindness

After .stcox for the multivariable model:
.stcurve, hazard at1(nblind=0) at2(nblind=1)yscale(log)



71

12.3b Summary of Multivariable Model Results

- Model estimates children with night-blind mothers have 41% increased adjusted hazard of death than children with mothers who are not night-blind.
- Conclusions?

The hazard "risk" of death is estimated to be 41% greater for an infant whose mother is night-blind (95% CI: 10 to 80% higher).