

Survival Analysis and Event History

Andy Grogan-Kaylor

11 Nov 2020

Introduction

“Survival analysis is a key technique in data-driven decision-making, which is now central to public interest because of COVID-19. Applying the correct technique for the specific question at hand is crucial for credible public health inferences. If you are interested in assessing how a risk factor or a potential treatment affects the progression of a disease—such as how long a patient takes to recover—then survival analysis techniques come into play. Survival analysis deeply respects the ultimate source of its data, often the disease experience or even the life and death of human patients. It seeks to exploit every last drop of information that this experience can render for saving lives—in particular, not only whether patients survived, but how long, and why. And it strives to do so with minimal assumptions, so that the data are truly driving the decision.”

—SAS Corporation

Key Concepts

WHO CARES how we measure time? Isn't it self-evident?

- Implementations differ; formulas are our friends
- $h(t) \sim x_1 + x_2 + \text{etc....}$: formula (effect on hazard (instantaneous rate of occurrence))

The “Hospital Bed Problem”

- Imagine a *Hypothetical Hospital*
- Imagine that there are 52 patients *total*.
- 51 of the patients are *long term patients*, who each stay for *1 year*.
- 1 of the patients is a *short term patient*, who stays for *1 week*.

Is this a hospital that serves mostly long-term, or short term patients?

```
. clear all

. set obs 52 // 52 hypothetical observations
number of observations (_N) was 0, now 52
```

```

. generate id = _n // set id = to observation #

. generate weeks = 52

. replace weeks = 1 if id == 52
(1 real change made)

. twoway (scatter id weeks if weeks == 52, msize(small)) /// staying 52 weeks
> (scatter id weeks if weeks == 1, msize(small)), /// staying 1 week
> title("Hypothetical Hospital") ///
> legend(on order(1 "long term" 2 "short term")) ///
> xtitle("week of discharge") ///
> ylabel(1(1)52, labels labsize(tiny) angle(horizontal) noticks nogrid) ///
> scheme(michigan)

. graph export hospital_bed_problem.png, width(1000) replace
(file hospital_bed_problem.png written in PNG format)

```

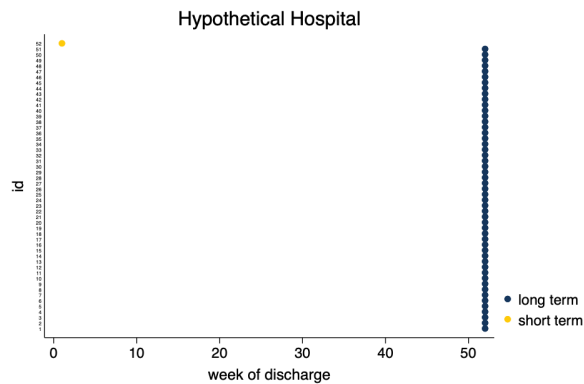


Figure 1: Illustration of Hospital Bed Problem

How To Measure Length of Stay (1)

```

. clear all

. set obs 25 // 25 hypothetical observations
number of observations (_N) was 0, now 25

. generate id = _n // set id = to observation #

. generate time = runiform(1, 100) // random times

. generate censored = time > 75 // censored if time > 75

. twoway (scatter id time if censored == 0) ///
> (scatter id time if censored == 1), ///
> title("Hypothetical Timing of Events") ///
> subtitle("Think About Different Kinds of Events") ///
> note("Study Ends At Time 75") ///
> legend(on order(1 "not censored" 2 "censored")) ///
> xline(75, lcolor("red")) /// censoring line at 75
> ylabel(1(1)25, labsize(vsmall) angle(horizontal)) /// lines from 1 to 25
> scheme(michigan)

. graph export timing_of_events.png, width(1000) replace
(file timing_of_events.png written in PNG format)

```

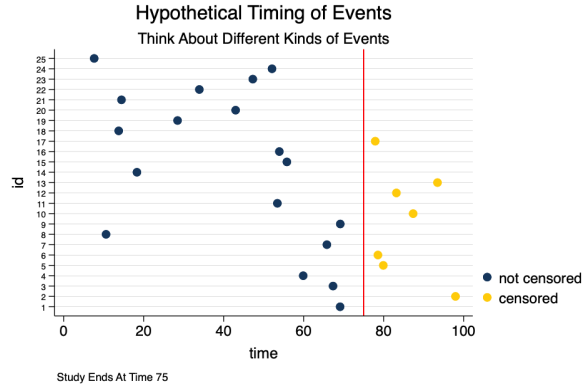


Figure 2: Timing Of Events

Animated

See [times-events-and-censoring.html](#)

How To Measure Length of Stay (2)

Event happened within a specified time (yes/no)

$$\ln\left(\frac{P(\text{event})}{1 - P(\text{event})}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e_i$$

- Statistically accurate, but we lose information on *when* the event happened.
- Statistically *less efficient*.

Time until Event

$$\text{time until event} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e_i$$

- What to do with events that haven't happened yet? (Censoring)
- Code as **missing**. Loss of information if using complete cases. Possible bias.
- Code as 0. Possible bias. They might happen at some point.
- Code as **time of censoring**. Possible bias. They might never happen. They might happen much later.

Hazard (Risk) of Event Occurrence

A more heuristic definition:

$$h(t) = \lim_{\delta \rightarrow 0} \frac{\text{probability of having an event before time } t + \delta}{\delta}$$

This definition per Johnson & Shih (2007)¹

¹Johnson, L. L., & Shih, J. H. (2007). CHAPTER 20 - An Introduction to Survival Analysis (J. I. Gallin & F. P. Ognibene, eds.). <https://doi.org/https://doi.org/10.1016/B978-012369440-9/50024-4>

A more formal definition:

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t | T > t)}{\Delta t}$$

This definition per Ragnar Frisch Centre for Economic Research (2020)²

A Policy Example (Welfare Reform, 1996)

From LaDonna Pavetti (1995)

- time in months
- new entrants (percent)
- all current recipients at a point in time (percent)

```
. clear all

. use Pavetti.dta
(Written by R.          )

. list, abbreviate(25) // list out the data
```

	time	new_entrants	all_current_recipients
1.	1-12	27.4	4.5
2.	13-24	14.8	4.8
3.	25-36	10	4.9
4.	37-48	7.7	5
5.	49-60	5.5	4.5
6.	Over 60	34.6	76.3

```
. graph bar (asis) all_current_recipients, /// this particular set of options was difficult to fig
> ure out!
> asyvars ///
> over(time) ///
> title("All Current Recipients") ///
> sub("By Months On Caseload") ///
> ytitle("percent") ///
> scheme(michigan)

. graph export all_current_recipients.png, width(1000) replace
(file all_current_recipients.png written in PNG format)
```

Welfare Reform (2)

```
. graph bar (asis) new_entrants, ///
> asyvars ///
> over(time) ///
> title("New Recipients") ///
> sub("By Months On Caseload") ///
> ytitle("percent") ///
> scheme(michigan)

. graph export new_recipients.png, width(1000) replace
(file new_recipients.png written in PNG format)
```

²Ragnar Frisch Centre for Economic Research (2020). Event History Analysis, Survival Analysis, Duration Analysis ,Transition Data Analysis, Hazard Rate Analysis. Oslo, Norway.

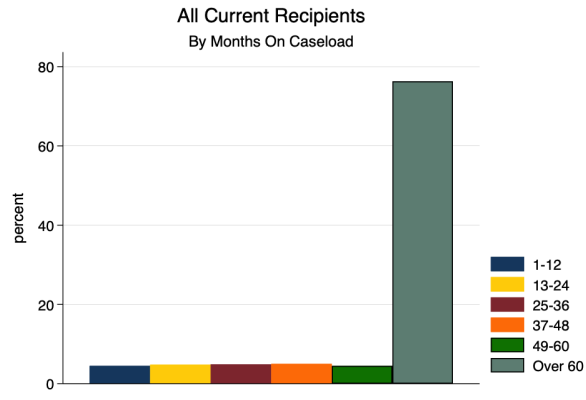


Figure 3: All Current Recipients by Months on Caseload

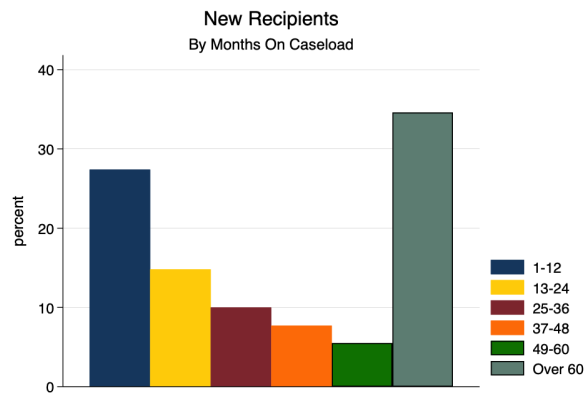


Figure 4: New Recipients by Months on Caseload

Musicians and Mortality (1)

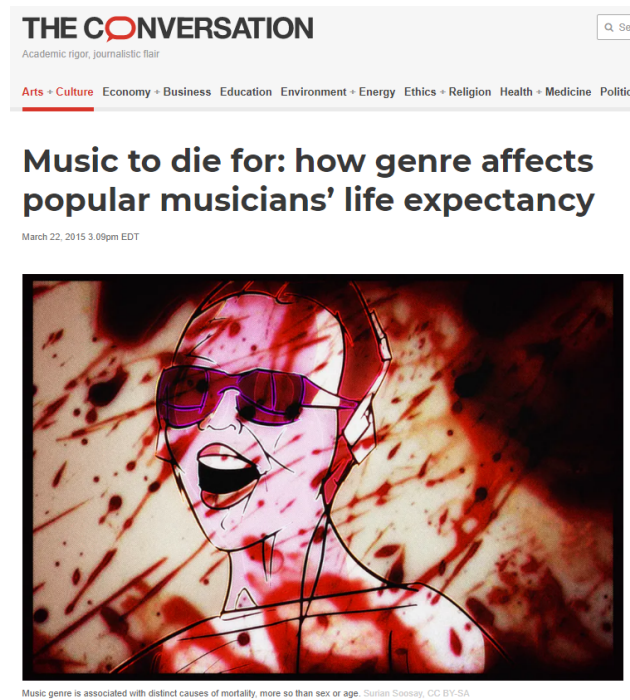


Figure 5: Music To Die For

Musicians and Mortality (2)

Cox Proportional Hazards Model

Formula

$h(t)$ the rate of occurrence.

$$h(t) = \lim_{\delta \rightarrow \infty} \frac{\text{probability of having an event before time } t + \delta}{\delta}$$

This definition per Johnson & Shih (2007)³

$$h(t) = h_0(t)e^{\beta_1 x_1 + \beta_2 x_2 + \text{etc.}}$$

We don't directly estimate the hazard, but estimate the effect of covariates on the hazard.

The event (birth, death, program entry, program departure) is coded as 1, so we are estimating the association of the covariates with event occurrence.

³Johnson, L. L., & Shih, J. H. (2007). CHAPTER 20 - An Introduction to Survival Analysis (J. I. Gallin & F. P. Ognibene, eds.). <https://doi.org/https://doi.org/10.1016/B978-012369440-9/50024-4>

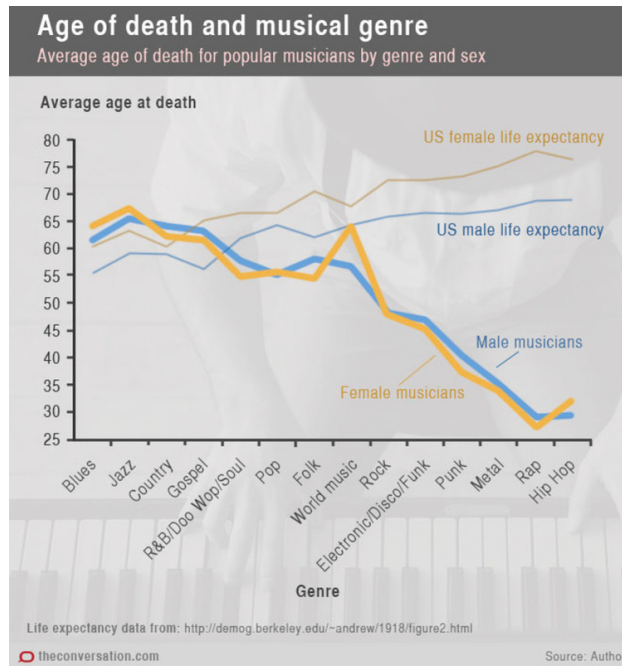


Figure 6: Musician Mortality

Cox Proportional Hazards Model in Stata

Using a data set referenced frequently in Stata `help` and Stata YouTube Videos

```
. clear all

. webuse drugtr // demonstration data set from Stata
(Patient Survival in Drug Trial)
```

Setup of Data

```
. stset // show st setup of data
-> stset studytime, failure(died)

    failure event:  died != 0 & died < .
obs. time interval:  (0, studytime]
exit on or before:  failure
```

```
48 total observations
0 exclusions
```

```
48 observations remaining, representing
31 failures in single-record/single-failure data
744 total analysis time at risk and under observation
      at risk from t =          0
earliest observed entry t =      0
last observed exit t =        39
```

Kaplan-Meier Survivor Function (per Gabriela Ortiz, Stata)

$$S(t) = Pr(T > t)$$

```
. sts graph, scheme(michigan) // Kaplan-Meier Survivor Function
    failure _d: died
    analysis time _t: studytime

. graph export survival0.png, width(1000) replace
(file survival0.png written in PNG format)
```

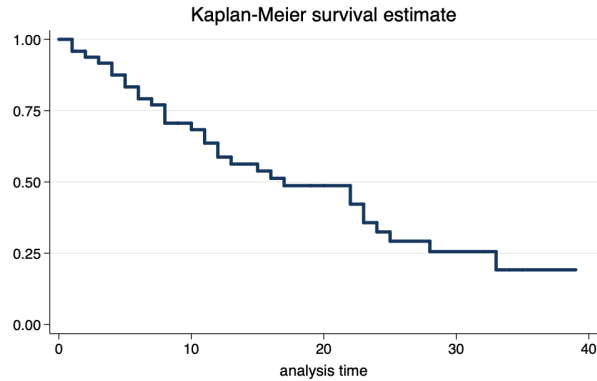


Figure 7: Kaplan-Meier Survivor Function

Cox Proportional Hazards Model

```
. stcox age drug // run Cox Proportional Hazards Model
    failure _d: died
    analysis time _t: studytime
Iteration 0:  log likelihood = -99.911448
Iteration 1:  log likelihood = -83.551879
Iteration 2:  log likelihood = -83.324009
Iteration 3:  log likelihood = -83.323546
Refining estimates:
Iteration 0:  log likelihood = -83.323546
Cox regression -- Breslow method for ties
No. of subjects =          48          Number of obs   =          48
No. of failures =          31
Time at risk   =          744
Log likelihood  =  -83.323546          LR chi2(2)       =          33.18
                                          Prob > chi2      =          0.0000
```

_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
age	1.120325	.0417711	3.05	0.002	1.041375	1.20526
drug	.1048772	.0477017	-4.96	0.000	.0430057	.2557622

Graph Survival Curves

```
. stcurve, survival scheme(michigan) // survival curve

. graph export survival1.png, width(1000) replace
(file survival1.png written in PNG format)

. stcurve, survival at1(drug=0) at2(drug=1) scheme(michigan) // survival curve by group

. graph export survival2.png, width(1000) replace
(file survival2.png written in PNG format)
```

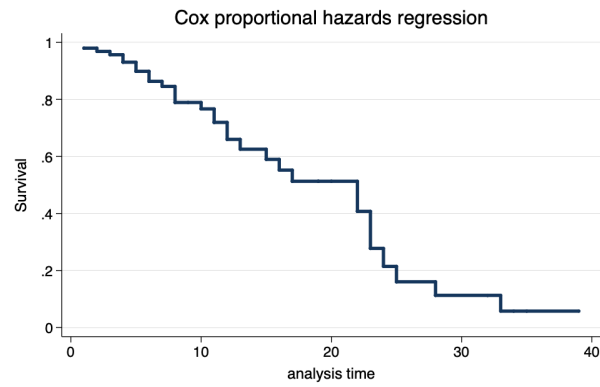



Figure 8: Survival Curve

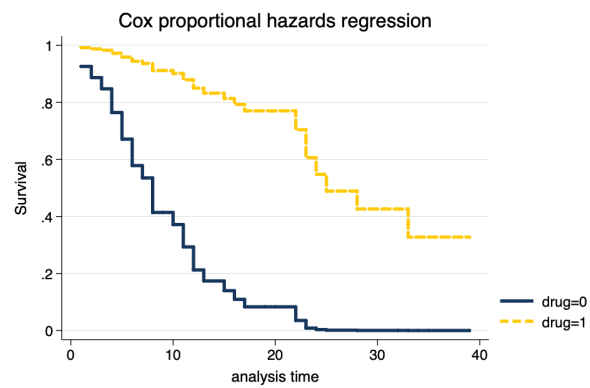


Figure 9: Survival Curve by Drug Group

Proportional Hazards Assumption

```
. estat phtest // formal test of PH assumption
```

Test of proportional-hazards assumption

Time: Time

	chi2	df	Prob>chi2
global test	0.43	2	0.8064

```
. sthplot, by(drug) scheme(michigan) // graphical test of PH assumption
```

failure _d: died

analysis time _t: studytime

```
. graph export ph.png, width(1000) replace
```

(file ph.png written in PNG format)

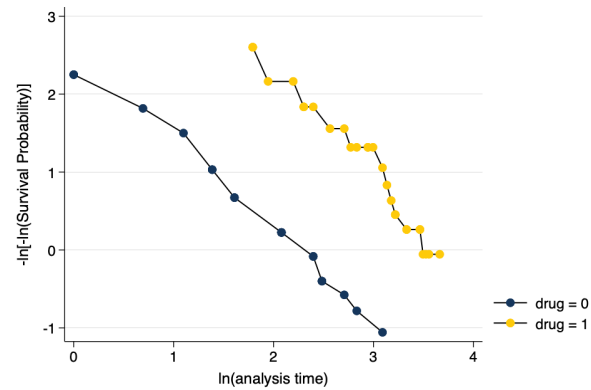


Figure 10: Graphical Assessment of Proportional Hazards Assumptions