

1 The problem of survival analysis

Survival analysis is concerned with analyzing the time to the occurrence of an event. For instance, we have a dataset in which the times are 1, 5, 9, 20, and 22. Perhaps those measurements are made in seconds, perhaps in days, but that does not matter. Perhaps the event is the time until a generator's bearings seize, the time until a cancer patient dies, or the time until a person finds employment, but that does not matter either.

For now, we will just abstract the underlying data-generating process and say that we have some times until an event occurs, and that those times are 1, 5, 9, 20, and 22. In addition, perhaps we have some covariates (additional variables) that we wish to use to "explain" these times. So, pretend that we have the following (completely made up) dataset:

time	x
1	3
5	2
9	4
20	9
22	10

Now, what is to keep us from simply analyzing these data using ordinary least-squares (OLS) linear regression? Why not simply fit the model

$$\text{time}_j = \beta_0 + \beta_1 x_j + \epsilon_j, \quad \epsilon_j \sim N(0, \sigma^2)$$

for $j = 1, \dots, 5$, or, alternatively,

$$\ln(\text{time}_j) = \beta_0 + \beta_1 x_j + \epsilon_j, \quad \epsilon_j \sim N(0, \sigma^2)$$

That is easy enough to do in Stata by typing

```
. regress time x
```

or

```
. generate lntime = ln(time)
. regress lntime x
```

These days, researchers would seldom analyze survival times in this manner, but why not? Before you answer too dismissively, we warn you that we, the authors, can think of problems for which this would be a perfectly reasonable model to use.

1.1 Parametric modeling

The problem with using OLS to analyze survival data lies with the assumed distribution of the residuals, ϵ_j . In linear regression, the residuals are assumed to be distributed normally, which is to say, time conditional on x_j is assumed to follow a normal distribution:

$$\text{time}_j \sim N(\beta_0 + \beta_1 x_j, \sigma^2), \quad j = 1, \dots, 5$$

The simple fact is that the assumed normality of time to an event is unreasonable for many events. It is unreasonable, for instance, if we are thinking about an event that has an instantaneous risk of occurring that is constant over time. In that case, the distribution of time would follow an exponential distribution. It is also unreasonable if we are analyzing survival times following a particularly serious surgical procedure. In that case, the distribution might have two modes: many patients die shortly after the surgery, but if they survive, the disease might be expected to return. One other problem is that a time to failure is always positive, while theoretically, the normal distribution is supported on the entire real line. Realistically, however, this fact alone is not enough to render the normal distribution useless in this context, since σ^2 may be chosen (or estimated) to make the probability of a negative failure time virtually zero.

At its core, survival analysis concerns nothing more than making a substitution for the normality assumption characterized by OLS with something more appropriate for the problem at hand.

Perhaps, if you were already familiar with survival analysis, when we asked “why not linear regression?”, you offered the excuse of right censoring—that in real data we often do not observe subjects long enough for all of them to fail. In our data, however, there was no censoring, and really, censoring is just a nuisance. We can fix linear regression easily enough to deal with right censoring. It goes under the name censored normal regression, and Stata’s `cnreg` command can fit such models; see [R] `tobit`. The real problem with linear regression in survival applications is with the assumed normality.

Not being already familiar with survival analysis, you might be tempted to use linear regression in the face of non-normality. Linear regression is known, after all, to be remarkably robust to deviations from normality, so why not just use it anyway? The problem is that the distributions for time to an event might be quite dissimilar from the normal—they are almost certainly nonsymmetric, they might be bimodal, and linear regression is not robust to these violations.

Substituting a more reasonable distributional assumption for ϵ_j leads to parametric survival analysis.

1.2 Semiparametric modeling

That results of analyses are being determined by the assumptions and not the data is always a source of concern, and this leads to a search for methods that do not require assumptions about the distribution of failure times. That, at first blush, seems hopeless. With survival data, the key insight into removing the distributional assumption is that, because events occur at given times, these events may be ordered and the analysis may be performed using the ordering of the survival times exclusively. Consider our dataset:

time	x
1	3
5	2
9	4
20	9
22	10

Examine the failure that occurred at time 1. Let's ask, "What is the probability of failure after exposure to the risk of failure for 1 unit of time?" At this point, observation 1 had failed, and the others had not. This reduces the problem to a problem of binary-outcome analysis,

time	x	outcome
1	3	1
5	2	0
9	4	0
20	9	0
22	10	0

and it would be perfectly reasonable for us to analyze failure at `time = 1` using, say, logistic regression

$$\begin{aligned}
 &= \Pr(\text{failure after exposure for 1 unit of time}) \\
 &= \Pr(\text{outcome}_j = 1) \\
 &= \frac{1}{1 + \exp(-\beta_0 - x_j\beta_x)}
 \end{aligned}$$

for $j = 1, \dots, 5$. This is easy enough to do:

```
. logistic outcome x
```

Do not make too much out of our choice of logistic regression—choose the analysis method you like. Use probit. Make a table. The point is that whatever particular technique you choose, you could do all your survival analysis using this analyze-the-first-failure method. It would be a mightily inefficient use of your data, but it would have the advantage that you would be making no assumptions about the distribution of failure times. Of course, you would have to give up on the idea of being able to make predictions conditional on x , but perhaps being able to predict whether failure occurs at `time = 1` would be sufficient.

There is nothing magical about the first death time; we could instead choose to analyze the second death time, which, it turns out in these data, is `time = 5`. We could

ask about the probability of failure, given exposure of 5 units of time, in which case we would exclude the first observation (which failed too early) and fit our logistic regression model using the second and subsequent observations:

```
. drop outcome
. generate outcome = cond(time==5,1,0) if time>=5
. logistic outcome x if time>=5
```

In fact, we could use this same procedure on each of the death times, separately.

Which analysis should we use? Well, there is slightly less information in the second analysis than in the first (because we have one less observation), and in the third than in the first two (for the same reason), and so on, so we should choose the first. It is, however, unfortunate, that we have to choose at all. Could we somehow combine all of these analyses and constrain the appropriate regression coefficients (say the coefficient on x) to be the same? The answer is yes, we could, and after some math, that leads to semiparametric survival analysis and, in particular, to Cox (1972) regression if a conditional logistic model is fit for each analysis. Conditional logistic models differ from ordinary logistic models for this example in that for the former we condition on the fact that we know that `outcome==1` for one and only one observation within each separate analysis.

However, for now we don't want to get lost in all the mathematical detail. What is important is that we could have done each of the analyses using whatever binary analysis method seemed appropriate. By doing so, we could combine them all if we are sufficiently clever in doing the math, and since each of the separate analyses made no assumption about the distribution of failure times, the combined analysis also makes no such assumption.

That last statement is rather slippery, so it does not hurt to verify its truth. We have been considering the data,

time	x
1	3
5	2
9	4
20	9
22	10

but now consider two variations on the data:

time	x
1.1	3
1.2	2
1.3	4
50.0	9
50.1	10

and

time	x
1	3
500	2
1000	4
10000	9
100000	10

These two alternatives have dramatically different distributions for time yet have the same temporal ordering and the same values of x . Think about performing the individual analyses on each of these datasets, and you will realize that the results you get will be exactly the same. Time plays no role other than ordering the observations.

The methods described above go under the name semiparametric analysis because, as far as time is concerned, they are nonparametric, but since we are still parameterizing the effect of x , there exists a parametric component to the analysis.

1.3 Nonparametric analysis

Semiparametric models are parametric in the sense that the effect of the covariates is still assumed to take a certain form. In the previous section, by performing a separate analysis at each failure time and concerning ourselves only with the order in which the failures occurred, we made no assumption about the distribution of time to failure. We did, however, make an assumption about how each subject's observed x value determined the probability that that subject would fail; for example, a probability determined by the logistic function.

An entirely nonparametric approach would be to do away with this assumption also and follow the philosophy of "letting the dataset speak for itself". There exists a vast literature on performing nonparametric regression using methods such as lowess or local polynomial regression; however, such methods do not adequately deal with censoring and other issues unique to survival data.

When no covariates exist, or when the covariates are qualitative in nature (gender, for instance), we can use nonparametric methods such as Kaplan and Meier (1958) or the method of Nelson (1972) and Aalen (1978) to estimate the probability of survival past a certain point in time, or to compare the survival experiences for each gender. These methods take into account censoring and other characteristics of survival data. There also exist methods such as the two-sample log-rank test, which can compare the survival experience across gender by using only the temporal ordering of the failure times. To wit, nonparametric methods make assumptions about neither (a) the distribution of the failure times nor (b) how covariates serve to change or shift the survival experience.

1.4 Linking the three approaches

Going back to our original data, consider the individual analyses we performed in order to obtain the semiparametric (combined) results. The individual analyses were

Pr(failure after exposure for (exactly) 1 unit of time)
 Pr(failure after exposure for (exactly) 5 units of time)
 Pr(failure after exposure for (exactly) 9 units of time)
 Pr(failure after exposure for (exactly) 20 units of time)
 Pr(failure after exposure for (exactly) 22 units of time)

We could omit any of the individual analyses above, and doing so would only affect the efficiency of our estimators. It is better, though, to include them all, so why not add the following to this list:

Pr(failure after exposure for (exactly) 1.1 units of time)
 Pr(failure after exposure for (exactly) 1.2 units of time)
 ...

That is, why not add individual analyses for all other times between the observed failure times? That would be a good idea because the more analyses we can combine, the more efficient our final results will be, which is to say that the standard errors of our estimated regression parameters will be smaller. The only reason we do not do this is that we do not know how to say anything about these intervening times—we do not know how to perform these analyses—unless we make an assumption about the distribution of failure time. If we made that assumption, we could perform the intervening analyses (the infinite number of them), and then we could combine them all to get super-efficient estimates. We could perform the individual analyses themselves a little differently, too, by taking into account the distributional assumptions, but that would only make our final analysis even more efficient.

That is the link between semiparametric and parametric analysis. Semiparametric analysis is nothing more than a combination of separate binary-outcome analyses, one per failure time, while parametric analysis is a combination of several analyses at *all* possible failure times. In parametric analysis, if no failures occur over a particular interval of time, that is informative. In semiparametric analysis, such periods are not informative. On the one hand, semiparametric analysis is advantageous in that it does not concern itself with the intervening analyses, yet parametric analysis will be more efficient if the proper distributional assumptions are made concerning those times when no failures are observed.

When no covariates are present, we hope that semiparametric methods such as Cox regression will produce estimates of relevant quantities (such as the probability of survival past a certain time) that are identical to the nonparametric estimates, and in fact, they do. When the covariates are qualitative in nature, parametric and semiparametric methods should yield more efficient tests and comparisons of the groups determined by the covariates than nonparametric methods, and these tests should agree. Should the tests disagree, this would serve as a signal that some of the assumptions made by the parametric or semiparametric models are incorrect.