

Parametric Survival Analysis

So far, we have focused primarily on nonparametric and semi-parametric approaches to survival analysis, with heavy emphasis on the Cox proportional hazards model:

$$\lambda(t, \mathbf{Z}) = \lambda_0(t) \exp(\beta \mathbf{Z})$$

We used the following estimating approach:

- We estimated $\lambda_0(t)$ nonparametrically, using the Kaplan-Meier estimator, or using the Kalbfleisch/Prentice estimator under the PH assumption
- We estimated β by assuming a linear model between the log HR and covariates, under the PH model

Both estimates were based on maximum likelihood theory.

Reading: for parametric models see Collett.

There are several reasons why we should consider some alternative approaches based on parametric models:

- The assumption of proportional hazards might not be appropriate (based on major departures)
- If a parametric model actually holds, then we would probably gain efficiency
- We may want to handle non-standard situations like
 - interval censoring
 - incorporating population mortality
- We may want to make some connections with other familiar approaches (e.g. use of the Poisson likelihood)
- We may want to obtain some estimates for use in designing a future survival study.

A simple start: Exponential Regression

- **Observed data:** $(X_i, \delta_i, \mathbf{Z}_i)$ for individual i ,
 $\mathbf{Z}_i = (Z_{i1}, Z_{i2}, \dots, Z_{ip})$ represents a set of p covariates.
- **Right censoring:** Assume that $X_i = \min(T_i, U_i)$
- **Survival distribution:** Assume T_i follows an exponential distribution with a parameter λ that depends on \mathbf{Z}_i , say $\lambda_i = \Psi(\mathbf{Z}_i)$. Then we can write:

$$T_i \sim \text{exponential}(\Psi(\mathbf{Z}_i))$$

First, let's review some facts about the exponential distribution (from our first survival lecture):

$$f(t) = \lambda e^{-\lambda t} \quad \text{for } t \geq 0$$

$$S(t) = P(T \geq t) = \int_t^{\infty} f(u) du = e^{-\lambda t}$$

$$F(t) = P(T < t) = 1 - e^{-\lambda t}$$

$$\lambda(t) = \frac{f(t)}{S(t)} = \lambda \quad \text{constant hazard!}$$

$$\Lambda(t) = \int_0^t \lambda(u) du = \int_0^t \lambda du = \lambda t$$

Now, we say that λ is a constant *over time* t , but we want to let it depend on the covariate values, so we are setting

$$\lambda_i = \Psi(\mathbf{Z}_i)$$

The hazard rate would therefore be the same for any two individuals with the same covariate values.

Although there are many possible choices for Ψ , one simple and natural choice is:

$$\Psi(\mathbf{Z}_i) = \exp[\beta_0 + Z_{i1}\beta_1 + Z_{i2}\beta_2 + \dots + Z_{ip}\beta_p]$$

WHY?

- ensures a positive hazard
- for an individual with $\mathbf{Z} = \mathbf{0}$, the hazard is e^{β_0} .

The model is called **exponential regression** because of the natural generalization from regular linear regression

Exponential regression for the 2-sample case:

- Assume we have only a single covariate $\mathbf{Z} = Z$, i.e., ($p = 1$).

Hazard Rate:

$$\Psi(\mathbf{Z}_i) = \exp(\beta_0 + Z_i\beta_1)$$

- Define:
 $Z_i = 0$ if individual i is in group 0
 $Z_i = 1$ if individual i is in group 1
- What is the hazard for group 0?
- What is the hazard for group 1?
- What is the hazard ratio of group 1 to group 0?
- What is the interpretation of β_1 ?

Likelihood for Exponential Model

Under the assumption of right censored data, each person has one of two possible contributions to the likelihood:

(a) they have an **event** at X_i ($\delta_i = 1$) \Rightarrow contribution is

$$L_i = \underbrace{S(X_i)}_{\text{survive to } X_i} \cdot \underbrace{\lambda(X_i)}_{\text{fail at } X_i} = e^{-\lambda X_i} \lambda$$

(b) they are **censored** at X_i ($\delta_i = 0$) \Rightarrow contribution is

$$L_i = \underbrace{S(X_i)}_{\text{survive to } X_i} = e^{-\lambda X_i}$$

The **likelihood** is the product over all of the individuals:

$$\begin{aligned}\mathcal{L} &= \prod_i L_i \\ &= \prod_i \underbrace{(\lambda e^{-\lambda X_i})^{\delta_i}}_{\text{events}} \underbrace{(e^{-\lambda X_i})^{(1-\delta_i)}}_{\text{censorings}} \\ &= \prod_i \lambda^{\delta_i} (e^{-\lambda X_i})\end{aligned}$$

Maximum Likelihood for Exponential

How do we use the likelihood?

- first take the log
- then take the partial derivative with respect to β
- then set to zero and solve for $\hat{\beta}$
- this gives us the **maximum likelihood estimators**

The log-likelihood is:

$$\begin{aligned}\log \mathcal{L} &= \log \left[\prod_i \lambda^{\delta_i} (e^{-\lambda X_i}) \right] \\ &= \sum_i [\delta_i \log(\lambda) - \lambda X_i] \\ &= \sum_i [\delta_i \log(\lambda)] - \sum_i \lambda X_i\end{aligned}$$

For the case of exponential regression, we now substitute the hazard $\lambda = \Psi(\mathbf{Z}_i)$ in the above log-likelihood:

$$\log \mathcal{L} = \sum_i [\delta_i \log(\Psi(\mathbf{Z}_i))] - \sum_i \Psi(\mathbf{Z}_i) X_i \quad (1)$$

General Form of Log-likelihood for Right Censored Data

In general, whenever we have right censored data, the likelihood and corresponding log likelihood will have the following forms:

$$\begin{aligned}\mathcal{L} &= \prod_i [\lambda_i(X_i)]^{\delta_i} S_i(X_i) \\ \log \mathcal{L} &= \sum_i [\delta_i \log(\lambda_i(X_i))] - \sum_i \Lambda_i(X_i)\end{aligned}$$

where

- $\lambda_i(X_i)$ is the hazard for the individual i who fails at X_i
- $\Lambda_i(X_i)$ is the cumulative hazard for an individual at their failure or censoring time

For example, see the derivation of the likelihood for a Cox model on p.11-18 of Lecture 4 notes. We started with the likelihood above, then substituted the specific forms for $\lambda(X_i)$ under the PH assumption.

Consider our model for the hazard rate:

$$\lambda = \Psi(\mathbf{Z}_i) = \exp[\beta_0 + Z_{i1}\beta_1 + Z_{i2}\beta_2 + \dots + Z_{ip}\beta_p]$$

We can write this using vector notation, as follows:

$$\begin{aligned} \text{Let } \mathbf{Z}_i &= (1, Z_{i1}, \dots, Z_{ip})^T \\ \text{and } \beta &= (\beta_0, \beta_1, \dots, \beta_p) \end{aligned}$$

(Since β_0 is the intercept (i.e., the log hazard rate for the baseline group), we put a “1” as the first term in the vector \mathbf{Z}_i .)

Then, we can write the hazard as:

$$\Psi(\mathbf{Z}_i) = \exp[\beta\mathbf{Z}_i]$$

Now we can substitute $\Psi(\mathbf{Z}_i) = \exp[\beta\mathbf{Z}_i]$ in the log-likelihood shown in (1):

$$\log \mathcal{L} = \sum_{i=1}^n \delta_i(\beta\mathbf{Z}_i) - \sum_{i=1}^n X_i \exp(\beta\mathbf{Z}_i)$$

Score Equations

Taking the derivative with respect to β_0 , the score equation is:

$$\frac{\partial \log \mathcal{L}}{\partial \beta_0} = \sum_{i=1}^n [\delta_i - X_i \exp(\beta \mathbf{Z}_i)]$$

For β_k , $k = 1, \dots, p$, the equations are:

$$\begin{aligned} \frac{\partial \log \mathcal{L}}{\partial \beta_k} &= \sum_{i=1}^n [\delta_i Z_{ik} - X_i Z_{ik} \exp(\beta \mathbf{Z}_i)] \\ &= \sum_{i=1}^n Z_{ik} [\delta_i - X_i \exp(\beta \mathbf{Z}_i)] \end{aligned}$$

To find the MLE's, we set the above equations to 0 and solve (simultaneously). The equations above imply that the MLE's are obtained by setting the weighted number of failures ($\sum_i Z_{ik} \delta_i$) equal to the weighted cumulative hazard ($\sum_i Z_{ik} \Lambda(X_i)$).

To find the variance of the MLE's, we need to take the second derivatives:

$$-\frac{\partial^2 \log \mathcal{L}}{\partial \beta_k \partial \beta_j} = \sum_{i=1}^n Z_{ik} Z_{ij} X_i \exp(\beta \mathbf{Z}_i)$$

Some algebra (see Cox and Oakes section 6.2) reveals that

$$\text{Var}(\hat{\beta}) = I(\beta)^{-1} = [\mathbf{Z}(\mathbf{I} - \Pi)\mathbf{Z}^T]^{-1}$$

where

- $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_n)$ is a $(p + 1) \times n$ matrix
(p covariates plus the “1” for the intercept β_0)
- $\Pi = \text{diag}(\pi_1, \dots, \pi_n)$ (this means that Π is a diagonal matrix, with the terms π_1, \dots, π_n on the diagonal)

- π_i is the probability that the i -th person is censored, so $(1 - \pi_i)$ is the probability that they failed.
- **Note:** The information $I(\beta)$ (inverse of the variance) is proportional to the number of failures, not the sample size. This will be important when we talk about study design.

The Single Sample Problem ($Z_i = 1$ for everyone):

First, what is the MLE of β_0 ?

We set $\frac{\partial \log \mathcal{L}}{\partial \beta_0} = \sum_{i=1}^n [\delta_i - X_i \exp(\beta_0 Z_i)]$ equal to 0 and solve:

$$\Rightarrow \sum_{i=1}^n \delta_i = \sum_{i=1}^n [X_i \exp(\beta_0)]$$

$$d = \exp(\beta_0) \sum_{i=1}^n X_i$$

$$\exp(\widehat{\beta_0}) = \frac{d}{\sum_{i=1}^n X_i}$$

$$\hat{\lambda} = \frac{d}{t}$$

where d is the total number of deaths (or events), and $t = \sum X_i$ is the total person-time contributed by all individuals.

If d/t is the MLE for λ , what does this imply about the MLE of β_0 ?

Using the previous formula $Var(\hat{\beta}) = [\mathbf{Z}(\mathbf{I} - \mathbf{\Pi})\mathbf{Z}^T]^{-1}$,
what is the variance of $\hat{\beta}_0$?:

With some matrix algebra, you can show that it is:

$$Var(\hat{\beta}_0) = \frac{1}{\sum_{i=1}^n (1 - \pi_i)} = \frac{1}{d}$$

What about $\hat{\lambda} = e^{\hat{\beta}_0}$?

By the delta method,

$$\begin{aligned} \text{Var}(\hat{\lambda}) &= \hat{\lambda}^2 \text{Var}(\hat{\beta}_0) \\ &= ? \end{aligned}$$

The Two-Sample Problem:

	Z_i	Subjects	Events	Follow-up
Group 0:	$Z_i = 0$	n_0	d_0	$t_0 = \sum_{i=1}^{n_0} X_i$
Group 1:	$Z_i = 1$	n_1	d_1	$t_1 = \sum_{i=1}^{n_1} X_i$

The log-likelihood:

$$\log \mathcal{L} = \sum_{i=1}^n \delta_i (\beta_0 + \beta_1 Z_i) - \sum_{i=1}^n X_i \exp(\beta_0 + \beta_1 Z_i)$$

$$\begin{aligned} \text{so } \frac{\partial \log \mathcal{L}}{\partial \beta_0} &= \sum_{i=1}^n [\delta_i - X_i \exp(\beta_0 + \beta_1 Z_i)] \\ &= (d_0 + d_1) - (t_0 e^{\beta_0} + t_1 e^{\beta_0 + \beta_1}) \end{aligned}$$

$$\begin{aligned} \frac{\partial \log \mathcal{L}}{\partial \beta_1} &= \sum_{i=1}^n Z_i [\delta_i - X_i \exp(\beta_0 + \beta_1 Z_i)] \\ &= d_1 - t_1 e^{\beta_0 + \beta_1} \end{aligned}$$

This implies: $\hat{\lambda}_1 = e^{\hat{\beta}_0 + \hat{\beta}_1} = ?$

$$\hat{\lambda}_0 = e^{\hat{\beta}_0} = ?$$

$$\hat{\beta}_0 = ?$$

$$\hat{\beta}_1 = ?$$

Important Result:

The maximum likelihood estimates (MLE's) of the hazard rates under the exponential model are the number of events divided by the person-years of follow-up!

(this result will be relied on heavily when we discuss study design)

Exponential Regression: Means and Medians

Mean Survival Time

For the exponential distribution, $E(T) = 1/\lambda$.

- **Control Group:**

$$\bar{T}_0 = 1/\hat{\lambda}_0 = 1/\exp(\hat{\beta}_0)$$

- **Treatment Group:**

$$\bar{T}_1 = 1/\hat{\lambda}_1 = 1/\exp(\hat{\beta}_0 + \hat{\beta}_1)$$

Median Survival Time

This is the value M at which $S(t) = e^{-\lambda t} = 0.5$, so
 $M = \text{median} = \frac{-\log(0.5)}{\lambda}$

- **Control Group:**

$$\hat{M}_0 = \frac{-\log(0.5)}{\hat{\lambda}_0} = \frac{-\log(0.5)}{\exp(\hat{\beta}_0)}$$

- **Treatment Group:**

$$\hat{M}_1 = \frac{-\log(0.5)}{\hat{\lambda}_1} = \frac{-\log(0.5)}{\exp(\hat{\beta}_0 + \hat{\beta}_1)}$$

Exponential Regression: Variance Estimates and Test Statistics

We can also calculate the variances of the MLE's as simple functions of the number of failures:

$$\text{var}(\hat{\beta}_0) = \frac{1}{d_0}$$

$$\text{var}(\hat{\beta}_1) = \frac{1}{d_0} + \frac{1}{d_1}$$

So our test statistics are formed as:

For testing $H_o : \beta_0 = 0$:

$$\begin{aligned}\chi_w^2 &= \frac{(\hat{\beta}_0)^2}{\text{var}(\hat{\beta}_0)} \\ &= \frac{[\log(d_0/t_0)]^2}{1/d_0}\end{aligned}$$

For testing $H_o : \beta_1 = 0$:

$$\begin{aligned}\chi_w^2 &= \frac{(\hat{\beta}_1)^2}{\text{var}(\hat{\beta}_1)} \\ &= \frac{\left[\log\left(\frac{d_1/t_1}{d_0/t_0}\right)\right]^2}{\frac{1}{d_0} + \frac{1}{d_1}}\end{aligned}$$

How would we form confidence intervals for the hazard ratio?

The Likelihood Ratio Test Statistic:

(An alternative to the Wald test)

A likelihood ratio test is based on 2 times the log of the ratio of the likelihoods under the null and alternative. We reject H_0 if $2 \log(\text{LR}) > \chi_{1,0.05}^2$, where

$$LR = \frac{\mathcal{L}(H_1)}{\mathcal{L}(H_0)} = \frac{\mathcal{L}(\hat{\lambda}_0, \hat{\lambda}_1)}{\mathcal{L}(\hat{\lambda})}$$

For a sample of n independent exponential random variables with parameter λ , the Likelihood is:

$$\begin{aligned} L &= \prod_{i=1}^n [\lambda^{\delta_i} \exp(-\lambda x_i)] \\ &= \lambda^d \exp(-\lambda \sum x_i) \\ &= \lambda^d \exp(-\lambda n \bar{x}) \end{aligned}$$

where d is the number of deaths or failures. The log-likelihood is

$$\ell = d \log(\lambda) - \lambda n \bar{x}$$

and the MLE is

$$\hat{\lambda} = d/(n\bar{x})$$

2-Sample Case: LR test calculations

Data:

Group 0: d_0 failures among the n_0 females
mean failure time is $\bar{x}_0 = (\sum_i^{n_0} X_i)/n_0$

Group 1: d_1 failures among the n_1 males
mean failure time is $\bar{x}_1 = (\sum_i^{n_1} X_i)/n_1$

Under the alternative hypothesis:

$$\mathcal{L} = \lambda_1^{d_1} \exp(-\lambda_1 n_1 \bar{x}_1) \times \lambda_0^{d_0} \exp(-\lambda_0 n_0 \bar{x}_0)$$

$$\log(\mathcal{L}) = d_1 \log(\lambda_1) - \lambda_1 n_1 \bar{x}_1 + d_0 \log(\lambda_0) - \lambda_0 n_0 \bar{x}_0$$

The MLE's are:

$$\hat{\lambda}_1 = d_1 / (n_1 \bar{x}_1) \quad \text{for males}$$

$$\hat{\lambda}_0 = d_0 / (n_0 \bar{x}_0) \quad \text{for females}$$

Under the null hypothesis:

$$\mathcal{L} = \lambda^{d_1 + d_0} \exp[-\lambda(n_1 \bar{x}_1 + n_0 \bar{x}_0)]$$

$$\log(\mathcal{L}) = (d_1 + d_0) \log(\lambda) - \lambda[n_1 \bar{x}_1 + n_0 \bar{x}_0]$$

The corresponding MLE is

$$\hat{\lambda} = (d_1 + d_0) / [n_1 \bar{x}_1 + n_0 \bar{x}_0]$$

A likelihood ratio test can be constructed by taking twice the difference of the log-likelihoods under the alternative and the null hypotheses:

$$-2 \left[(d_0 + d_1) \log \left(\frac{d_0 + d_1}{t_0 + t_1} \right) - d_1 \log[d_1/t_1] - d_0 \log[d_0/t_0] \right]$$

Nursing home example:

For the females:

- $n_0 = 1173$
- $d_0 = 902$
- $t_0 = 310754$
- $\bar{x}_0 = 265$

For the males:

- $n_1 = 418$
- $d_1 = 367$
- $t_1 = 75457$
- $\bar{x}_1 = 181$

Plugging these values in, we get a LR test statistic of 64.20.

Hand Calculations using events and follow-up:

By adding up “LOS” for males to get t_1 and for females to get t_0 , I obtained:

- $d_0 = 902$ (females)
 $d_1 = 367$ (males)
- $t_0 = 310754$ (female follow-up)
 $t_1 = 75457$ (male follow-up)
- This yields an estimated log HR:

$$\hat{\beta}_1 = \log \left[\frac{d_1/t_1}{d_0/t_0} \right] = \log \left[\frac{367/75457}{902/310754} \right] = \log(1.6756) = 0.5162$$

- The estimated standard error is:

$$\sqrt{\text{var}(\hat{\beta}_1)} = \sqrt{\frac{1}{d_1} + \frac{1}{d_0}} = \sqrt{\frac{1}{902} + \frac{1}{367}} = 0.06192$$

- So the Wald test becomes:

$$\chi_W^2 = \frac{\hat{\beta}_1^2}{\text{var}(\hat{\beta}_1)} = \frac{(0.51619)^2}{0.061915} = 69.51$$

- We can also calculate $\hat{\beta}_0 = \log(d_0/t_0) = -5.842$,
along with its standard error $se(\hat{\beta}_0) = \sqrt{(1/d_0)} = 0.0333$

Exponential Regression in STATA

```
. use nurshome
. stset los fail
. streg gender, dist(exp) nohr
      failure _d:  fail
      analysis time _t:  los
Iteration 0:  log likelihood = -3352.5765
Iteration 1:  log likelihood = -3321.966
Iteration 2:  log likelihood = -3320.4792
Iteration 3:  log likelihood = -3320.4766
Iteration 4:  log likelihood = -3320.4766
Exponential regression -- log relative-hazard form
No. of subjects =          1591          Number of obs   =          1591
No. of failures =          1269
Time at risk    =          386211
Log likelihood  =  -3320.4766          LR chi2(1)       =          64.20
                                          Prob > chi2     =          0.0000
-----
```

_t	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
gender	.516186	.0619148	8.337	0.000	.3948352	.6375368
_cons	-5.842142	.0332964	-175.459	0.000	-5.907402	-5.776883

```
-----
```

Since $Z = 8.337$, the chi-square test is $Z^2 = 69.51$.

The Weibull Regression Model

At the beginning of the course, we saw that the survivorship function for a Weibull random variable is:

$$S(t) = \exp[-\lambda(t^\kappa)]$$

and the hazard function is:

$$\lambda(t) = \kappa \lambda t^{(\kappa-1)}$$

The Weibull regression model assumes that for someone with covariates \mathbf{Z}_i , the survivorship function is

$$S(t; \mathbf{Z}_i) = \exp[-\Psi(\mathbf{Z}_i)(t^\kappa)]$$

where $\Psi(\mathbf{Z}_i)$ is defined as in exponential regression to be:

$$\Psi(\mathbf{Z}_i) = \exp[\beta_0 + Z_{i1}\beta_1 + Z_{i2}\beta_2 + \dots Z_{ip}\beta_p]$$

For the 2-sample problem, we have:

$$\Psi(\mathbf{Z}_i) = \exp[\beta_0 + Z_{i1}\beta_1]$$

Weibull MLEs for the 2-sample problem:

Log-likelihood:

$$\log \mathcal{L} = \sum_{i=1}^n \delta_i \log [\kappa \exp(\beta_0 + \beta_1 Z_i) X_i^{\kappa-1}] - \sum_{i=1}^n X_i^{\kappa} \exp(\beta_0 + \beta_1 Z_i)$$

$$\Rightarrow \exp(\hat{\beta}_0) = d_0/t_0\kappa \quad \exp(\hat{\beta}_0 + \hat{\beta}_1) = d_1/t_1\kappa$$

where

$$t_{j\kappa} = \sum_{i=1}^{n_j} X_i^{\hat{\kappa}} \quad \text{among } n_j \text{ subjects}$$

$$\hat{\lambda}_0(t) = \hat{\kappa} \exp(\hat{\beta}_0) t^{\hat{\kappa}-1} \quad \hat{\lambda}_1(t) = \hat{\kappa} \exp(\hat{\beta}_0 + \hat{\beta}_1) t^{\hat{\kappa}-1}$$

$$\begin{aligned} \widehat{HR} &= \hat{\lambda}_1(t)/\hat{\lambda}_0(t) = \exp(\hat{\beta}_1) \\ &= \exp\left(\frac{d_1/t_1\kappa}{d_0/t_0\kappa}\right) \end{aligned}$$

Weibull Regression: Means and Medians

Mean Survival Time

For the Weibull distribution, $E(T) = \lambda^{(-1/\kappa)} \Gamma[(1/\kappa) + 1]$.

- **Control Group:**

$$\bar{T}_0 = \hat{\lambda}_0^{(-1/\hat{\kappa})} \Gamma[(1/\hat{\kappa}) + 1]$$

- **Treatment Group:**

$$\bar{T}_1 = \hat{\lambda}_1^{(-1/\hat{\kappa})} \Gamma[(1/\hat{\kappa}) + 1]$$

Median Survival Time

For the Weibull distribution, $M = \text{median} = \left[\frac{-\log(0.5)}{\lambda} \right]^{1/\kappa}$

- **Control Group:**

$$\hat{M}_0 = \left[\frac{-\log(0.5)}{\hat{\lambda}_0} \right]^{1/\hat{\kappa}}$$

- **Treatment Group:**

$$\hat{M}_1 = \left[\frac{-\log(0.5)}{\hat{\lambda}_1} \right]^{1/\hat{\kappa}}$$

where $\hat{\lambda}_0 = \exp(\hat{\beta}_0)$ and $\hat{\lambda}_1 = \exp(\hat{\beta}_0 + \hat{\beta}_1)$.

Note: the symbol Γ is the “gamma” function. If x is an integer, then

$$\Gamma(x) = (x - 1)!$$

In cases where x is not an integer, this function has to be evaluated numerically. In homework and labs, I will supply this value to you.

The Weibull regression model is very easy to fit:

- In STATA: Just specify `dist(weibull)` instead of `dist(exp)` within the `streg` command
- In SAS: use model option `dist=weibull` within the `proc lifereg` procedure

Note: to get more information on these modeling procedures, use the online help facilities. For example, in STATA, you can type:

```
.help streg
```


Weibull in Stata:

```
. streg gender, dist(weibull) nohr
```

```
        failure _d:  fail  
analysis time _t:  los
```

Fitting constant-only model:

```
Iteration 0:  log likelihood = -3352.5765  
Iteration 1:  log likelihood = -3074.978  
Iteration 2:  log likelihood = -3066.1526  
Iteration 3:  log likelihood = -3066.143  
Iteration 4:  log likelihood = -3066.143
```

Fitting full model:

```
Iteration 0:  log likelihood = -3066.143  
Iteration 1:  log likelihood = -3045.8152  
Iteration 2:  log likelihood = -3045.2772  
Iteration 3:  log likelihood = -3045.2768  
Iteration 4:  log likelihood = -3045.2768
```

Weibull regression -- log relative-hazard form

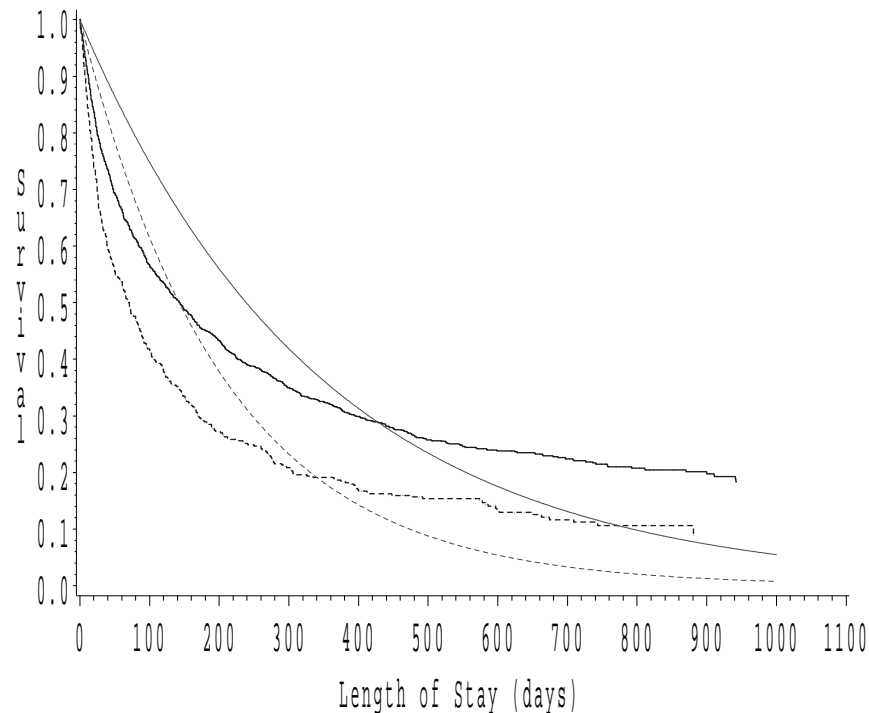
No. of subjects = 1591 Number of obs = 1591
 No. of failures = 1269
 Time at risk = 386211
 LR chi2(1) = 41.73
 Log likelihood = -3045.2768 Prob > chi2 = 0.0000

_t	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
gender	.4138082	.0621021	6.663	0.000	.2920903	.5355261
_cons	-3.536982	.0891809	-39.661	0.000	-3.711773	-3.362191
/ln_p	-.4870456	.0232089	-20.985	0.00	-.5325343	-.4415569
p	.614439	.0142605			.5871152	.6430345
1/p	1.627501	.0377726			1.555127	1.703243

Comparison of Exponential with Kaplan-Meier

We can see how well the Exponential model fits by comparing the survival estimates for males and females under the exponential model, i.e., $P(T \geq t) = e^{(-\hat{\lambda}_z t)}$, to the Kaplan-Meier survival estimates:

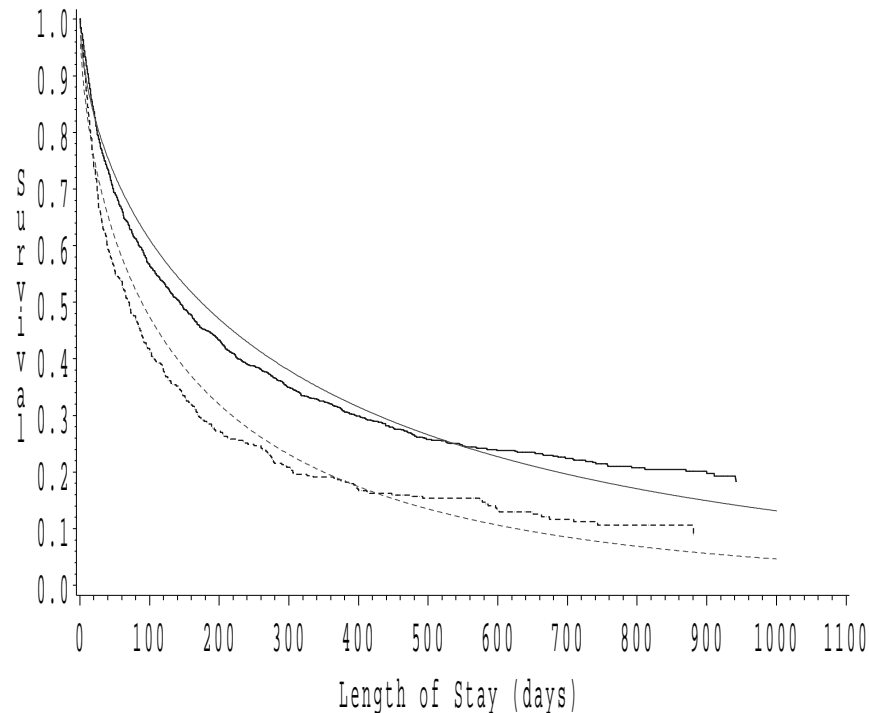
Predicted Survival for Exponential model vs Kaplan-Meier



Comparison of Weibull with Kaplan-Meier

We can see how well the Weibull model fits by comparing the survival estimates, $P(T \geq t) = e^{(-\hat{\lambda}_z t^{\hat{\kappa}})}$, to the Kaplan-Meier survival estimates.

Predicted Survival for Weibull model vs Kaplan-Meier



Which do you think fits best?

Other useful plots for evaluating fit to exponential and Weibull models

- $-\log(\hat{S}(t))$ vs t
- $\log[-\log(\hat{S}(t))]$ vs $\log(t)$

Why are these useful?

If T is exponential, then $S(t) = \exp(-\lambda t)$

$$\text{so } \log(S(t)) = -\lambda t$$

$$\text{and } \Lambda(t) = \lambda t$$

a straight line in t with slope λ and intercept=0

If T is Weibull, then $S(t) = \exp(-(\lambda t)^\kappa)$

$$\text{so } \log(S(t)) = -\lambda t^\kappa$$

$$\text{then } \Lambda(t) = \lambda t^\kappa$$

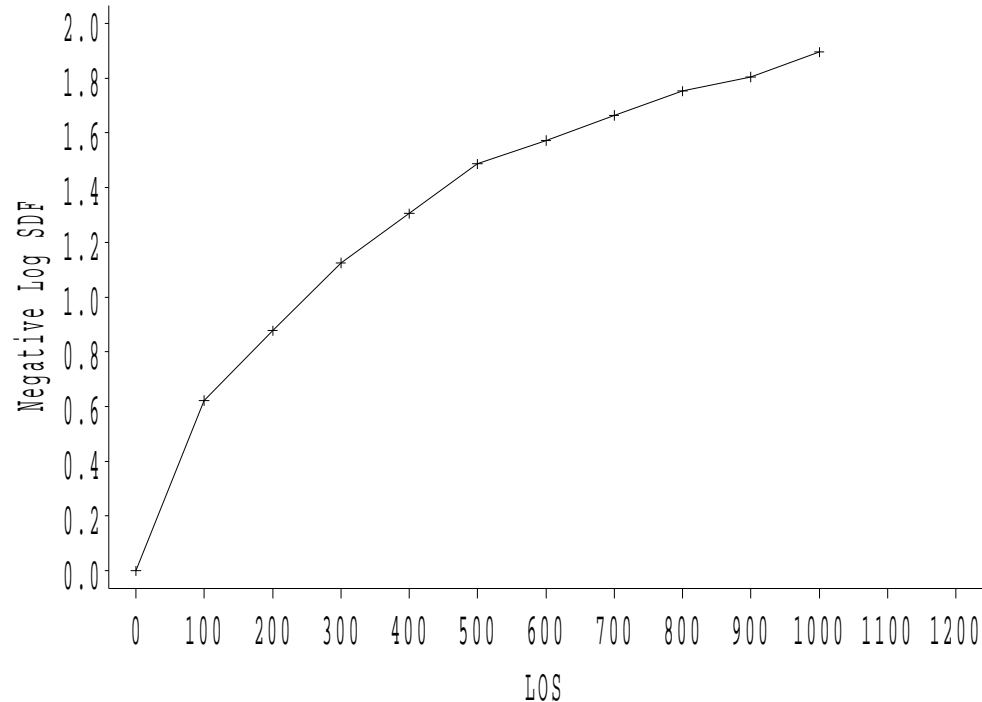
$$\text{and } \log(-\log(S(t))) = \log(\lambda) + \kappa * \log(t)$$

a straight line in $\log(t)$ with slope κ and intercept $\log(\lambda)$.

So we can calculate our estimated $\Lambda(t)$ and plot it versus t , and if it seems to form a straight line, then the exponential distribution is probably appropriate for our dataset.

Plots for nursing home data: $\hat{\Lambda}(t)$ vs t

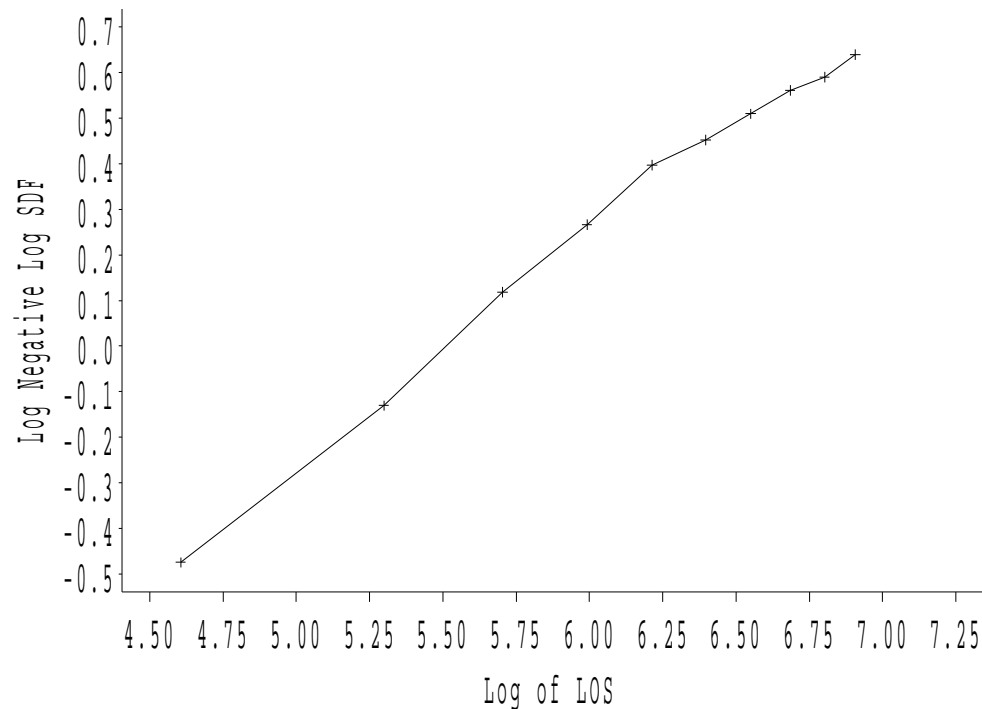
Estimated cumulative hazard vs time



Or we can plot $\log \hat{\Lambda}(t)$ versus $\log(t)$, and if it seems to form a straight line, then the Weibull distribution is probably appropriate for our dataset.

Plots for nursing home data: $\log[-\log(\hat{S}(t))]$ vs $\log(t)$

Estimated log cumulative hazard vs log time



Comparison of Methods for the Two-sample problem:

Data:

	Z_i	Subjects	Events	Follow-up
Group 0:	$Z_i = 0$	n_0	d_0	$t_0 = \sum_{i=1}^{n_0} X_i$
Group 1:	$Z_i = 1$	n_1	d_1	$t_1 = \sum_{i=1}^{n_1} X_i$

In General:

$$\lambda_z(t) = \lambda(t, Z = z) \quad \text{for } z = 0 \text{ or } 1.$$

The hazard rate depends on the value of the covariate Z . In this case, we are assuming that we only have a single covariate, and it is binary ($Z = 1$ or $Z = 0$)

Reading from Collett:

Section(s)	Description
4.1.1, 4.1.2	Exponential properties
4.1.3	Weibull properties
4.3.1, 4.4.2	Exponential ML estimation
4.3.2	Weibull ML estimation
4.5	General Weibull regression
4.6	Model selection - Weibull regression
4.7	Weibull/AFT model connection
Ch.6	AFT - Other parametric models

MODELS

Exponential Regression:

$$\lambda_z(t) = \exp(\beta_0 + \beta_1 Z)$$

$$\Rightarrow \lambda_0 = \exp(\beta_0)$$

$$\lambda_1 = \exp(\beta_0 + \beta_1)$$

$$HR = \exp(\beta_1)$$

Weibull Regression:

$$\lambda_z(t) = \kappa \exp(\beta_0 + \beta_1 Z) t^{\kappa-1}$$

$$\Rightarrow \lambda_0 = \kappa \exp(\beta_0) t^{\kappa-1}$$

$$\lambda_1 = \kappa \exp(\beta_0 + \beta_1) t^{\kappa-1}$$

$$HR = \exp(\beta_1)$$

Proportional Hazards Model:

$$\lambda_z(t) = \lambda_0(t) \exp(\beta_1)$$

$$\Rightarrow \lambda_0 = \lambda_0(t) \quad \text{KM?}$$

$$\lambda_1 = \lambda_0(t) \exp(\beta_1)$$

$$HR = \exp(\beta_1)$$

Remarks

- Exponential model is a special case of the Weibull model with $\kappa = 1$ (note: Collett uses γ instead of κ)
- Exponential and Weibull models are both special cases of the Cox PH model.

How can you show this?

- If either the exponential model or the Weibull model is valid, then these models will tend to be more efficient than PH (smaller s.e.'s of estimates). This is because they assume a particular form for $\lambda_0(t)$, rather than estimating it at every death time.

For the Exponential model, the hazards are constant over time, given the value of the covariate Z_i :

$$Z_i = 0 \Rightarrow \hat{\lambda}_0 = \exp(\hat{\beta}_0)$$

$$Z_i = 1 \Rightarrow \hat{\lambda}_0 = \exp(\hat{\beta}_0 + \hat{\beta}_1)$$

For the Weibull model, we have to estimate the hazard as a function of time, given the estimates of β_0, β_1 and κ :

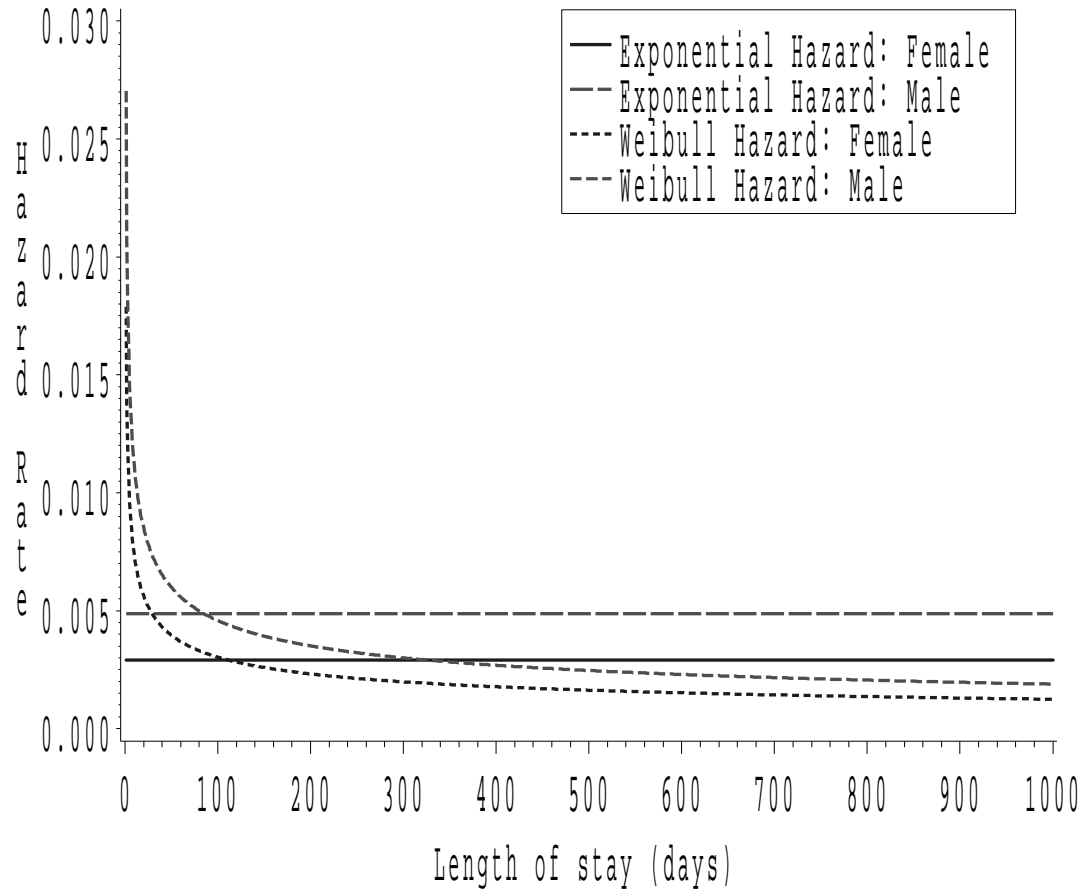
$$Z_i = 0 \Rightarrow \hat{\lambda}_0(t) = \hat{\kappa} \exp(\hat{\beta}_0) t^{\hat{\kappa}-1}$$

$$Z_i = 1 \Rightarrow \hat{\lambda}_1(t) = \hat{\kappa} \exp(\hat{\beta}_0 + \hat{\beta}_1) t^{\hat{\kappa}-1}$$

However, the ratio of the hazards is still just $\exp(\hat{\beta}_1)$, since the other terms cancel out.

Here's what the estimated hazards look like for the nursing home data:

Estimated Hazards for Weibull & Exponential by Gender



Proportional Hazards Model:

To get the MLE's for this model, we have to maximize the Cox partial likelihood iteratively. There are not closed form expressions like above.

$$\begin{aligned} L(\beta) &= \prod_{i=1}^n \left[\frac{e^{\beta \mathbf{Z}_i}}{\sum_{\ell \in \mathcal{R}(X_i)} e^{\beta \mathbf{Z}_\ell}} \right]^{\delta_i} \\ &= \prod_{i=1}^n \left[\frac{e^{\beta_0 + \beta_1 Z_i}}{\sum_{\ell \in \mathcal{R}(X_i)} e^{\beta_0 + \beta_1 Z_\ell}} \right]^{\delta_i} \end{aligned}$$

Comparison with Proportional Hazards Model

```
. stcox gender, nohr
      failure _d:  fail
      analysis time _t:  los
Iteration 0:  log likelihood = -8556.5713
Iteration 1:  log likelihood = -8537.8013
Iteration 2:  log likelihood = -8537.5605
Iteration 3:  log likelihood = -8537.5604
Refining estimates:
Iteration 0:  log likelihood = -8537.5604
Cox regression -- Breslow method for ties
No. of subjects =          1591          Number of obs   =          1591
No. of failures =          1269
Time at risk    =          386211
Log likelihood  =  -8537.5604          LR chi2(1)       =          38.02
                                          Prob > chi2     =          0.0000
-----
```

	<u>_t</u>					
	<u>_d</u>	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
gender		.3943588	.0621004	6.350	0.000	.2726441 .5160734

```
-----
```

For the PH model, $\hat{\beta}_1 = 0.394$ and $\widehat{HR} = e^{0.394} = 1.483$.

Comparison with the Logrank and Wilcoxon Tests

```
. sts test gender
      failure _d: fail
      analysis time _t: los
```

Log-rank test for equality of survivor functions

```
-----
```

	Events	
gender	observed	expected
0	902	995.40
1	367	273.60
Total	1269	1269.00

```
-----
```

chi2(1) = 41.08
Pr>chi2 = 0.0000

```
. sts test gender, wilcoxon
```

```
      failure _d: fail  
analysis time _t: los
```

```
Wilcoxon (Breslow) test for equality of survivor functions
```

```
-----  
      | Events                Sum of  
gender | observed      expected      ranks  
-----+-----  
0      |      902          995.40      -99257  
1      |      367          273.60       99257  
-----+-----  
Total  |      1269         1269.00           0
```

```
      chi2(1) =      41.47
```

```
      Pr>chi2 =      0.0000
```

Comparison of Hazard Ratios and Test Statistics for effect of Gender

Model/Method	λ_0	λ_1	HR	log(HR)	se(log HR)	Wald Statistic
Exponential	0.0029	0.0049	1.676	0.5162	0.0619	69.507
Weibull						
$t = 50$	0.0040	0.0060	1.513	0.4138	0.0636	42.381
$t = 100$	0.0030	0.0046	1.513			
$t = 500$	0.0016	0.0025	1.513			
Logrank						41.085
Wilcoxon						41.468
Cox PH						
Ties=Breslow			1.483	0.3944	0.0621	40.327
Ties=Discrete			1.487	0.3969	0.0623	40.565
Ties=Efron			1.486	0.3958	0.0621	40.616
Ties=Exact			1.486	0.3958	0.0621	40.617
Score (Discrete)						41.085

Comparison of Mean and Median Survival Times by Gender

Model/Method	Mean Survival		Median Survival	
	Female	Male	Female	Male
Exponential	344.5	205.6	238.8	142.5
Weibull	461.6	235.4	174.2	88.8
Kaplan-Meier	318.6	200.7	144	70
Cox PH (Kalbfleisch/Prentice)			131	72

The Accelerated Failure Time Model

The general form of an accelerated failure time (AFT) model is:

$$\log(T_i) = \beta_{AFT} \mathbf{Z}_i + \sigma \epsilon$$

where

- $\log(T_i)$ is the log of a survival time
- β_{AFT} is the vector of AFT model parameters corresponding to the covariate vector \mathbf{Z}_i
- ϵ is a random “error” term
- σ is a scale factor

In other words, we can model the log-survival times as a linear function of the covariates!

The `streg` command in STATA (without the `exponential` or `weibull` option) uses this “log-linear” model for fitting parametric models.

By choosing different distributions for ϵ , we can obtain different parametric distributions:

- Exponential
- Weibull
- Gamma
- Log-logistic
- Normal
- Lognormal

We can compare the predicted survival under any of these parametric distributions to the KM estimated survival to see which one seems to fit best.

Once we decide on a certain class of model (say, Gamma), we can evaluate the contributions of covariates by finding the MLE's, and constructing Wald, Score, or LR tests of the covariate effects.

We can motivate the AFT model by first demonstrating the following two relationships:

1. For the Exponential Model:

If the failure times $T_i = T(\mathbf{Z}_i)$ follow an exponential distribution, i.e., $S_i(t) = e^{-\lambda_i t}$ with $\lambda_i = \exp(\beta \mathbf{Z}_i)$, then

$$\log(T_i) = -\beta \mathbf{Z}_i + \epsilon$$

where ϵ follows an extreme value distribution (which just means that e^ϵ follows a unit exponential distribution).

2. For the Weibull Model:

If the failure times $T_i = T(\mathbf{Z}_i)$ follow a Weibull distribution, i.e., $S_i(t) = e^{-\lambda_i t^\kappa}$ with $\lambda_i = \exp(\beta \mathbf{Z}_i)$, then

$$\log(T_i) = -\sigma \beta \mathbf{Z}_i + \sigma \epsilon$$

where ϵ again follows an extreme value distribution, and $\sigma = 1/\kappa$.

In other words, both the Exponential and Weibull model can be written in the form of a log-linear model for the survival times, if we choose the right distribution for ϵ .

The log-linear form for the exponential can be derived by:

- (1) Creating a new variable $T_0 = T_Z \times \exp(\beta \mathbf{Z}_i)$
- (2) Taking the log of T_Z , yielding $\log(T_Z) = \log\left(\frac{T_0}{\exp(\beta \mathbf{Z}_i)}\right)$

Step (1): For an exponential model, recall that:

$$S_i(t) = Pr(T_Z \geq t) = e^{-\lambda t}, \quad \text{with } \lambda = \exp(\beta \mathbf{Z}_i)$$

It follows that $T_0 \sim \text{exp}(1)$:

$$\begin{aligned} S_0(t) = Pr(T_0 \geq t) &= Pr(T_Z \cdot \exp(\beta \mathbf{Z}) \geq t) \\ &= Pr(T_Z \geq t \exp(-\beta \mathbf{Z})) \\ &= \exp[-\lambda t \exp(-\beta \mathbf{Z})] \\ &= \exp[-\exp(\beta \mathbf{Z}) t \exp(-\beta \mathbf{Z})] \\ &= \exp(-t) \end{aligned}$$

Step (2): Now take the log of the survival time:

$$\begin{aligned}\log(T_Z) &= \log\left(\frac{T_0}{\exp(\beta\mathbf{Z}_i)}\right) \\ &= \log(T_0) - \log(\exp(\beta\mathbf{Z}_i)) \\ &= -\beta\mathbf{Z}_i + \log(T_0) \\ &= -\beta\mathbf{Z}_i + \epsilon\end{aligned}$$

where $\epsilon = \log(T_0)$ follows the **extreme value** distribution.

Relationship between Exponential and Weibull

If T_Z has a Weibull distribution, i.e., $S(t) = e^{-\lambda t^\kappa}$ with $\lambda = \exp(\beta \mathbf{Z}_i)$, then you can show that the new variable

$$T_Z^* = T_Z^\kappa$$

follows an exponential distribution with parameter $\exp(\beta \mathbf{Z}_i)$. Based on the previous page, we can therefore write:

$$\log(T^*) = -\beta \mathbf{Z} + \epsilon$$

(where ϵ has an extreme value distribution.)

But since $\log(T^*) = \log(T^\kappa) = \kappa \times \log(T)$, we can write:

$$\begin{aligned}\log(T) &= \log(T^*)/\kappa \\ &= (1/\kappa) (-\beta \mathbf{Z}_i + \epsilon) \\ &= -\sigma \beta \mathbf{Z}_i + \sigma \epsilon\end{aligned}$$

where $\sigma = 1/\kappa$.

This motivates the following general definition of the **Accelerated Failure Time Model** by:

$$\log(T_i) = \beta_{AFT} \mathbf{Z}_i + \sigma \epsilon$$

where ϵ is a random “error” term, σ is a scale factor, Y is the log of a survival random variable, and

$$\beta_{AFT} = -\sigma \beta_e$$

where β_e came from the hazard $\lambda = \exp(\beta \mathbf{Z})$.

The defining feature of an AFT model is:

$$S(t; \mathbf{Z}) = S_i(t) = S_0(\phi t)$$

That is, the effect of covariates is to accelerate (stretch) or decelerate (shrink) the time-scale.

Effect of AFT on hazard:

$$\lambda_i(t) = \phi \lambda_0(\phi t)$$

One way to interpret the AFT model is via its effect on median survival times. If $S_i(t) = 0.5$, then $S_0(\phi t) = 0.5$. This means:

$$M_i = \phi M_0$$

Interpretation:

- For $\phi < 1$, there is an acceleration of the endpoint (if $M_0 = 2$ yrs in control and $\phi = 0.5$, then $M_i = 1$ yr.)
- For $\phi > 1$, there is a stretching or delay in endpoint
- In general, the lifetime of individual i is ϕ times what they would have experienced in the reference group

Since ϕ must be positive and a function of the covariates, we model $\phi = \exp(\beta \mathbf{Z}_i)$.

When does Proportional hazards = AFT?

According to the proportional hazards model:

$$S(t) = S_0(t)^{\exp(\beta \mathbf{Z}_i)}$$

and according to the accelerated failure time model:

$$S(t) = S_0(t \exp(\beta \mathbf{Z}_i))$$

Say $T_i \sim Weibull(\lambda, \kappa)$. Then $\lambda(t) = \lambda \kappa t^{(\kappa-1)}$

Under the AFT model:

$$\begin{aligned}\lambda_i(t) &= \phi \lambda_0(\phi t) \\ &= e^{\beta \mathbf{Z}_i} \lambda_0(e^{\beta \mathbf{Z}_i} t) \\ &= e^{\beta \mathbf{Z}_i} \lambda_0 \kappa (e^{\beta \mathbf{Z}_i} t)^{(\kappa-1)} \\ &= (e^{\beta \mathbf{Z}_i})^\kappa \lambda_0 \kappa t^{(\kappa-1)} \\ &= (e^{\beta \mathbf{Z}_i})^\kappa \lambda_0(t)\end{aligned}$$

But this looks just like the PH model:

$$\lambda_i(t) = \exp(\beta^* \mathbf{Z}_i) \lambda_0(t)$$

It turns out that the Weibull distribution (and exponential, since this is just a special case of a Weibull with $\kappa = 1$) is the only one for which the accelerated failure time and proportional hazards models coincide.

Special cases of AFT models

- Exponential regression: $\sigma = 1$, ϵ following the extreme value distribution.
- Weibull regression: σ arbitrary, ϵ following the extreme value distribution.
- Lognormal regression: σ arbitrary, ϵ following the normal distribution.

Examples in stata: Using the STREG command, one has the following options of distributions for the log-survival times:

```
. streg trt, dist(lognormal)
```

- exponential
- weibull
- gompertz
- lognormal
- loglogistic
- gamma

```
. streg gender, dist(exponential) nohr
```

_t	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
gender	.516186	.0619148	8.337	0.000	.3948352	.6375368

```
. streg gender, dist(weibull) nohr
```

_t	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
gender	.4138082	.0621021	6.663	0.000	.2920903	.5355261
1/p	1.627501	.0377726			1.555127	1.703243

```
. streg gender, dist(lognormal)
```

_t	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
gender	-.6743434	.1127352	-5.982	0.000	-.8953002	-.4533866
_cons	4.957636	.0588939	84.179	0.000	4.842206	5.073066
sigma	1.94718	.040584			1.86924	2.028371

```
. streg gender, dist(gamma)
```

_t	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
gender	-.6508469	.1147116	-5.674	0.000	-.8756774	-.4260163
_cons	4.788114	.1020906	46.901	0.000	4.58802	4.988208
sigma	1.97998	.0429379			1.897586	2.065951

This gives a good idea of the sensitivity of the test of gender to the choice of model. It is also easy to get predicted survival curves under any of the parametric models using the following:

```
. streg gender, dist(gamma)
. stcurv, survival
```

The options HAZARD and CUMHAZ can also be substituted for SURVIVAL above to obtain plots.