

# Survival analysis and censored data \*

- we observe either the survival time  $T$  or else the censoring time  $C \Rightarrow Y = \min(T, C)$
- we also observe a status indicator  $\delta = \begin{cases} 1 & \text{if } T < C \\ 0 & \text{if } T > C \end{cases}$  then if  $\delta = 1$  we observe the true survival time, if  $\delta = 0$ , we observe the censoring time.
- we have  $n$  pairs  $(Y_i, \delta_i) \dots (Y_n, \delta_n)$

In general, we need to assume that the censoring mechanism is independent: conditional on the features, the event time  $T$  is independent of the censoring time  $C$ . we often have to consider the data collection in order to determine whether independent censoring is a reasonable assumption.

## The Kaplan-Meier Survival curve

the survival curve (a function) is defined as  $S(t) = P(T > t)$ .  
 we let  $d_1 < d_2 < \dots < d_k$  denote the  $k$  unique death times among the non censored patients and we let  $q_h$  denote the number of patients who died at time  $d_h$ . For  $h = 1, \dots, k$  we let  $n_h$  denote the number of patients alive and in the study just before  $d_h$ ; these are the at risk patients. The set of patients that are at risk at a given time are referred to as the risk set

$$P(T > d_h) = P(T > d_h | T > d_{h-1}) P(T > d_{h-1}) + P(T > d_h | T \leq d_{h-1}) P(T \leq d_{h-1})$$

$$S(d_h) = P(T > d_h) = P(T > d_h | T > d_{h-1}) \times \dots \times P(T > d_2 | T > d_1) P(T > d_1) \text{ with } P(T > d_j | T > d_{j-1}) = \frac{n_j - q_j}{n_j}$$

$$\hat{S}(t) = \prod_{j=1}^k \left( \frac{n_j - q_j}{n_j} \right) \text{ and for times } t \text{ between } d_h \text{ and } d_{h+1}, \text{ we set } \hat{S}(t) = \hat{S}(d_h)$$

## The log-rank test

we examine how the events in each group unfold sequentially in time. In order to test  $H_0: E(X) = \mu$  for some random variable  $X$ , one approach is to construct a test statistic of the form  $w = \frac{\bar{X} - \mu}{\sqrt{\text{Var}(X)}}$

	group 1	group 2	Total
died	$q_{1h}$	$q_{2h}$	$q_h$
survival	$n_{1h} - q_{1h}$	$n_{2h} - q_{2h}$	$n_h - q_h$
total	$n_{1h}$	$n_{2h}$	$n_h$

$$\Rightarrow w = \sum_{h=1}^k \left( q_{1h} - \frac{q_h}{n_h} n_{1h} \right)$$

$$x = \sum_{h=1}^k q_{1h} \text{ is } \mu = \sum_{h=1}^k \frac{q_{1h}}{n_h} q_h$$

$$\frac{1}{\sum_{h=1}^k} \frac{q_h (n_{1h}/n_h) (1 - n_{1h}/n_h) (n_h - q_h)}{n_{h-1}}$$

when the sample size is large,  $w$  has a standard normal distribution. This can be used to compute a p-value for the null hypothesis that there is no difference between the survival curves in the two groups.

## Regression model with survival curve

The hazard function (a rate) is defined as  $\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t < T \leq t + \Delta t | T > t)}{\Delta t}$  where  $T$  is the (uncensored) survival time.  
 $\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t < T \leq t + \Delta t) \wedge (T > t) / \Delta t}{P(T > t)} = \frac{P(t < T \leq t + \Delta t) / \Delta t}{P(T > t)} = \frac{f(t)}{S(t)}$   
 where  $f(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t < T \leq t + \Delta t)}{\Delta t}$  is the probability density function associated with  $T$ , i.e. the instantaneous rate of death at time  $t$ .

The likelihood associated with the  $i$ th observation is

$L_i = \delta(t_i)$  if the  $i$ th observation is not censored

$S(t_i)$  if the  $i$ th observation is censored

$= \delta(t_i)^{\delta_i} S(t_i)^{1-\delta_i}$

The likelihood is if  $Y = t_i$  and the observation is not censored then the likelihood is the probability of dying in a time interval around time  $t_i$ . If the  $i$ th observation is censored, then the likelihood is the probability of surviving at least until time  $t_i$ .

To model the survival time as function of the covariates  $\Rightarrow h(t|x_i) = \exp(\beta_0 + \sum_{j=1}^p \beta_j x_{ij})$

$\Rightarrow h(t|x_i) = h_0(t) \exp(\sum_{j=1}^p x_{ij} \beta_j) \Rightarrow$  one unit increase in  $x_{ij}$  corresponds to an increase in  $h(t|x_i)$  by a factor of  $\exp(\beta_j)$

The partial likelihood  $PL(\beta) = \prod_{i: \delta_i=1} \frac{\exp(\sum_{j=1}^p x_{ij} \beta_j)}{\sum_{i: t_i \geq t_i} \exp(\sum_{j=1}^p x_{ij} \beta_j)}$  To estimate  $\beta$ , we simply maximize PL with respect to  $\beta$ .

Shrinkage for the Cox model

$-\log \left( \prod_{i: \delta_i=1} \frac{\exp(\sum_{j=1}^p x_{ij} \beta_j)}{\sum_{i: t_i \geq t_i} \exp(\sum_{j=1}^p x_{ij} \beta_j)} \right) + \lambda P(\beta)$  is minimized with respect to  $\beta = (\beta_1, \dots, \beta_p)^T$ , if  $P(\beta) = \sum_{j=1}^p \beta_j^2 = \text{ridge}$ ,  $\sum_{j=1}^p |\beta_j| = \text{lasso}$ .