

# Lecture 1 Basic Concepts

①

Study "time to event"  $T$ , where event

- age @ death
- age @ recurrence
- cure (time to cure after treatment)
- time to pregnancy after fertility treatment
- time to equipment failure
- ...

Characteristics of 'survival data'

- $T \geq 0$
- often 'incomplete' data due to censoring

Note - not the same as 'missing data' (bad measurement, - here 'missingness' provides partial information about the outcome)

Typical data, 2 treatment groups - control - treatment (receive screening for breast cancer)  
follow for 18 years

During this time, some women die from breast cancer  
Some are still alive after 18 years

(1) So ... for women who died from BC we know the exact time to event. For women still alive we only know  $T > 18$  years.

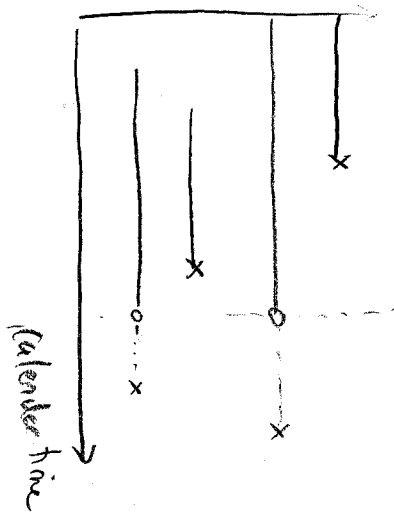
Possible complications

- censoring during study due to other (random) causes
  - loss to follow-up (perhaps moved)
- censoring due to death due to other causes
  - car accident - unrelated
  - heart disease - perhaps related  $\Rightarrow$  competing risks

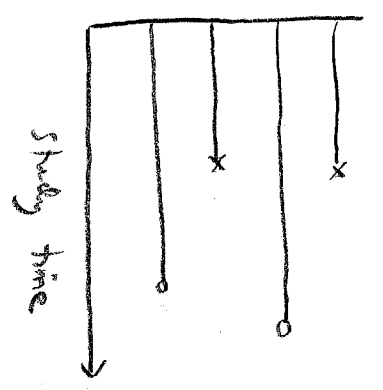
- different length of study for different women  
 $\rightarrow$  perhaps recruitment to study took ~ 5 years.  
Then some women have been followed for only 13 years, some for more.

Important to define start and end of study properly.

18 years



$\Rightarrow$



(more later)

# Goals

③

- ① Estimate & interpret survival characteristics
- ② Compare survival between groups
- ③ Assess relationship between survival and various predictors.

## Assumptions

- Independent observations
- Independent censoring

Why a separate course?

Typical summary statistics, like mean, are biased when data are censored.

⇒ General approach: estimate distribution of  $T$  from censored data and base estimates of the mean etc on this

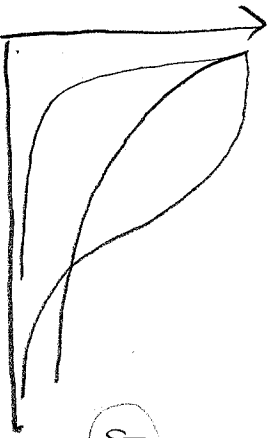
## Describing T

•  $S(t) = P(T > t)$

- survival function

-  $= 1 - F(t) = 1 - P(T \leq t)$

-  $S(0) = 1$ ,  $S(\infty) = 0$  (can adjust this to  $S(0) > 0$  if there is a cure prob.)



discuss

For continuous T  $\Rightarrow$  Definition  $F(t) = \int_0^t f(u) du$

(11)

$$\Rightarrow S(t) = \int_t^{\infty} f(u) du \Rightarrow f(t) = -\frac{dS(t)}{dt}$$

- Slope of survival fun.

Sample data  $\Rightarrow F_n(t) = \sum_{i=1}^n \frac{1}{n} 1\{t_i \leq t\}$

$$S_n(t) = \sum_{i=1}^n \frac{1}{n} 1\{t_i > t\}$$

= % (pop) of sample survival times that exceeds  $t$ .

$$\rightarrow S(t) \text{ as } n \rightarrow \infty$$

Other quantities of interest

mean survival time  $E(T) = \int_0^{\infty} t f(t) dt$

$$= - \int_0^{\infty} t dS(t) dt$$

integr by parts

$$\int_0^{\infty} S(t) dt$$

median survival time  $t_M: S(t_M) = \frac{1}{2}$

o mean residual life time (mrl)

(5)

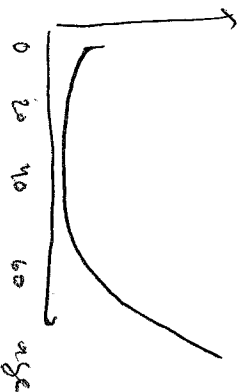
$$mrl(t_0) = E[T - t_0 | T \geq t_0]$$

$$\begin{aligned} &= \frac{\int_{t_0}^{\infty} (u - t_0) f(u) du}{S(t_0)} = - \frac{\int_{t_0}^{\infty} (u - t_0) dS(u)}{S(t_0)} \\ &\approx \frac{\int_{t_0}^{\infty} S(u) du}{S(t_0)} \end{aligned}$$

Hazard rate

discrete version - mortality rate

$$m(t) = P[E \leq T < E+1 | T \geq t]$$



human mortality rates

$$\text{hazard rate } h(t) = \lim_{\Delta t \rightarrow 0} \frac{P[t \leq T < t + \Delta t | T \geq t]}{\Delta t}$$

so for small  $\Delta t$   $P[t \leq T < t + \Delta t | T \geq t] \approx h(t) \cdot \Delta t$

$$\text{Note } h(t) = \lim_{\Delta t \rightarrow 0} \frac{P[t \leq T < t + \Delta t]}{\Delta t} = \frac{f(t)}{S(t)} = - \frac{S'(t)}{S(t)}$$

$$= - \frac{d \log S(t)}{dt}$$

## Cumulative hazard, $H(t)$

(6)

$$\text{hazard rate } h(t) = - \frac{d \log S(t)}{dt}$$

$$\Rightarrow H(t) = \int_0^t h(u) du = - \log \{S(t)\}$$

$$\Rightarrow S(t) = e^{-H(t)} = e^{-\int_0^t h(u) du}$$

[ Discrete case:

Definition of hazard rate

$$h(t) = \frac{P(t)}{S(t)} = P(T=t | T \geq t)$$

$$\begin{aligned} \text{Note } \circ \text{ We can write } \frac{S(t_j)}{S(t_{j-1})} &= \frac{P(T > t_j)}{P(T > t_{j-1})} = \frac{P(T > t_j | T > t_{j-1})}{P(T > t_{j-1})} \\ &= P(T > t_j | T > t_{j-1}) \end{aligned}$$

$$\circ \text{ Follows that } S(t) = \prod_{t_j \leq t} \frac{S(t_j)}{S(t_{j-1})} = \prod_{t_j \leq t} S(t_j | T > t_{j-1})$$

$$= \prod_{t_j \leq t} \{1 - h(t_j)\}$$

$$\text{(imp } S(t) = e^{-\int_0^t h(u) du})$$

$$\text{(Cumulative hazard } H(t) = \sum_{t_j \leq t} h(t_j) \text{ (KM)}$$

or alt  $\approx$  for small  $h$

$$\sum_{t_j \leq t} \ln \{1 - h(t_j)\} \text{ (CSD)}$$

$\Rightarrow S(t) = e^{-H(t)}$  from here as well

Remarks:  $M(t) \leftrightarrow S(t)$

hazard rate  $\neq$  a probability - more like a density

Some examples

o Constant hazard  $M(t) = \lambda$

$$\Rightarrow S(t) = e^{-\int_0^t \lambda du} = e^{-\lambda t}$$

$$f(t) = \frac{-dS(t)}{dt} = \lambda e^{-\lambda t}$$

The exponential distribution!

For Exp. dist  $\Rightarrow E(T) = \frac{1}{\lambda}$

$$M(t) = \int_{t_0}^{\infty} \frac{e^{-\lambda t} dt}{e^{-\lambda t_0}} = \frac{1}{\lambda} = E(T)$$

$$\Rightarrow P(T > t+z | T \geq t) = \frac{P(T > t+z)}{P(T \geq t)} = \frac{e^{-\lambda(t+z)}}{e^{-\lambda t}} = e^{-\lambda z} = P(T > z)$$

Memoryless

In practice, this distribution is too skewed most of the time.

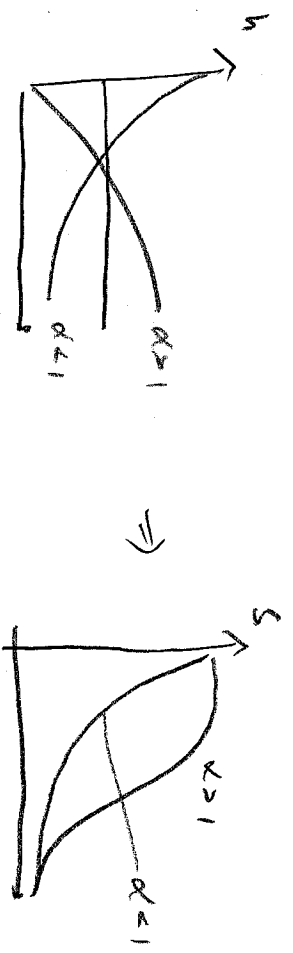
Reality check; plot  $t$  vs  $\log \tilde{S}(t)$   $\rightarrow$  should be a straight line

Weibull  $h(t) = \alpha \lambda t^{\alpha-1}$ ,  $S(t) = e^{-\lambda t^\alpha}$

$\alpha = 1 \Rightarrow$  weverts back to exponential

$\alpha > 1 \Rightarrow h(t)$  increasing  $\Rightarrow$  'aging'

$\alpha < 1 \Rightarrow h(t)$  decreasing (less common)

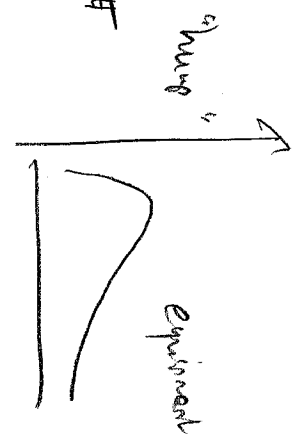


Note, if Weibull has  $\rightarrow \log(-\log S(t)) = \log \lambda + \alpha \log t$   
 Verify check; plot  $\ln t$  vs  $\ln(-\ln S(t))$

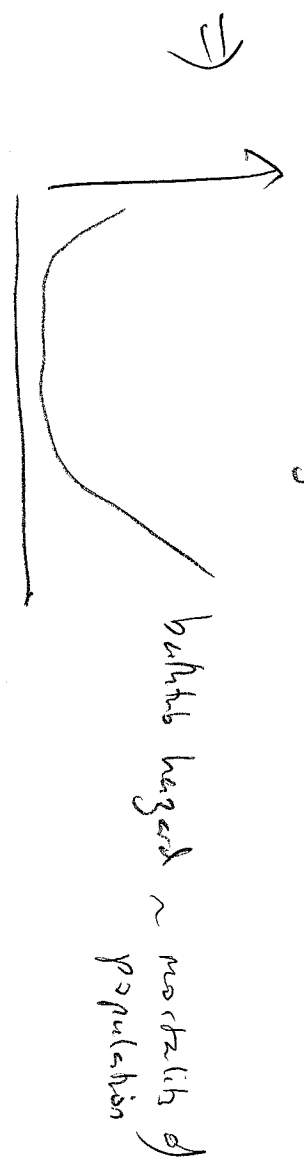
Other parametric forms

log-normal  $S(t) = 1 - \Phi \left[ \frac{\ln t - \mu}{\sigma} \right]$

log-logistic  $S(t) = \frac{e^{-(y-\frac{y_0}{\sigma})}}{1 + e^{-(y-\frac{y_0}{\sigma})}}$ ,  $y = \ln t$



Other commonly observed hazards





Looking ahead a little — regression models

(9)

①  $T \geq 0$  so standard regression is inappropriate  
 $\Rightarrow$  log transformation (accelerated failure time model)

$$Y = \ln T = \mu + X\beta + \varepsilon$$

$\underbrace{\hspace{2cm}}$  model form (parametric)  $\swarrow$  if specified  $\rightarrow$  parametric  
 $\searrow$  if not  $\rightarrow$  semi-parametric

Note;  $P(T > t | X) = ?$  quantity of interest

$$\bullet P(T > t | X = 0) = S_0(t) = P(e^{(\mu + \varepsilon)} > t)$$

$$\bullet P(T > t | X) = P(Y > \ln t | X) = P(Y - X\beta > \ln t - X\beta | X)$$

$$= P(e^{Y - X\beta} > t e^{-X\beta} | X) = P(e^{(\mu + \varepsilon)} > t e^{-X\beta} | X)$$

$$= \int_0^{\infty} t e^{-X\beta} \}$$

Same survival function, but time-scale has been compressed / elongated (Dep on sign of  $\beta$ )

$$\text{If } \beta > 0, e^{-X\beta} < 1, t e^{-X\beta} < t \text{ life has been 'compressed'}$$

$$\text{If } \beta < 0, \text{ there has been 'stretching' by covariate}$$

Estimation of  $\beta$  models can be a bit tricky  
so the more popular regression model for survival is  $\Rightarrow$

## ② Proportional hazard

(10)

$$h(t|X) = h_0(t) c(X\beta) \quad \text{where } c(X\beta) \text{ st. } h \text{ positive}$$

$$\text{Usually } c(X\beta) = e^{X\beta}$$

$$\text{Note } \frac{h(t|X_1)}{h(t|X_2)} = \exp((X_1 - X_2)\beta) \quad \text{not a function of } t$$

$$\text{Example } X = \begin{matrix} 0 & \text{control} \\ 1 & \text{treatment} \end{matrix}$$

$$\neq \boxed{\text{relative risk RR}} \\ = \frac{h_0(t)e^\beta}{h_0(t)} = e^\beta$$

Constant across time

$$\eta \quad \beta > 0 \Rightarrow \text{RR of group } X=1 > 1$$

$$\eta \quad \beta < 0 \Rightarrow \text{RR of group } X=1 < 1$$

(treatment effective)

- Estimation (as we shall see) — can estimate  $\beta$ 's w/o knowing  $h_0(t) \Rightarrow$  then use  $\hat{\beta}$  and  $\hat{S}(t)$  to estimate  $h_0(t)$ .

$$\begin{aligned} \bullet S(t|X) &= \exp\left(-\int_0^t h_0(u) e^{X\beta} du\right) = \exp\left(-e^{X\beta} \int_0^t h_0(u) du\right) = \exp(-e^{X\beta} H_0(t)) \\ &= \exp(-e^{X\beta} \ln S_0(t)) = \exp\left(-\ln S_0(t) e^{X\beta}\right) = S_0(t) e^{X\beta} \end{aligned}$$

- Proportional hazard = assumption  $\Rightarrow$  Need to check!  
(more in Regression later)