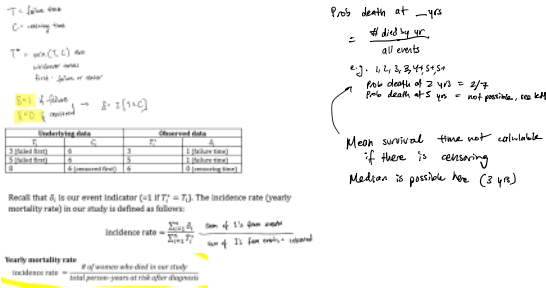


**Event:** The outcome of interest  
**Time origin:** the beginning of the survival time ("time zero")  
**Censoring:** A subject is censored when the endpoint of interest has not been observed for the individual  
**Right censoring:** Censor at the last time the participant was observed  
**Non-informative/independent censoring:** At each censoring time, those who are censored have the same prognosis as those who remain under observation  
**Administrative censoring:** When a study ends and the participant did not experience the event by the time of study end. A type of non-informative censoring.  
**Random censoring:** Censoring times are unassociated with T  
**Informative/dependent censoring:** Those left under observation have systematically different failure times than those who are censored. Kaplan-Meier fails if this is true.

**Rate-based analyses:**  
Related to simple survival analysis methods. Cons: Assumes the event rate is CONSTANT.



**Survival Function**

T = A survival random variable that measures the time elapsed from an origin ("time zero") until the event of interest.  
T is positive (T > 0). T can either be discrete or continuous.

**S(t)** = Survival function aka survival curve.  
 $S(t) = \Pr(T > t)$ . Probability that an individual will "survive" beyond a given length of time t.  
Key properties:  
•  $S(t) = \Pr(T > t) = 1 - \Pr(T \leq t) = 1 - F(t)$   
•  $0 \leq S(t) \leq 1$  for all t. Since S(t) is a probability.  
•  $S(0) = 1$  (100% survival). Since everyone is still at-risk at time 0,  
•  $S(\infty) = 0$ . (0% survival after infinite follow-up.) Assuming everyone eventually fails.  
• S(t) decreases or stays constant over time but never increases. If t < u, you can survive to time u only if you survive to time t, so  $S(t) \geq S(u)$

**F(t)** = cumulative distribution function (CDF) of T.  
 $F(t) = \Pr(T \leq t)$ . Probability that an individual will get outcome before or on time t.

For n uncensored individuals:  
$$\hat{F}(t) = \frac{1}{n} \sum_{j: t_j \leq t} 1 \longrightarrow \text{POINT ESTIMATE e.g. } \frac{13 \text{ deaths}}{21 \text{ ppl total}}$$

**Estimations:**  
 $E[\hat{F}(t)] = F(t)$   
 $\text{Var}(\hat{F}(t)) = \frac{1}{n} F(t)(1 - F(t))$

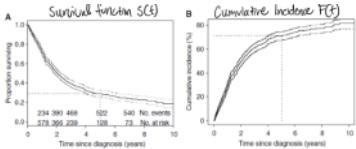
**Wald 95% CI:**  
$$\hat{F}(t) \pm 1.96 \sqrt{\frac{1}{n} \hat{F}(t)(1 - \hat{F}(t))}$$

**Log-transformed Wald 95% CI:**  
$$\log(\hat{F}(t)) \pm 1.96 \sqrt{\text{Var}(\log(\hat{F}(t)))}$$
  
$$= \hat{F}(t) \exp\left[\pm 1.96 \sqrt{\text{Var}(\log(\hat{F}(t)))}\right]$$

**Log-log transformed Wald 95% CI:**  
$$\log(-\log(\hat{F}(t))) \pm 1.96 \sqrt{\text{Var}(\log(-\log(\hat{F}(t))))}$$
  
$$= [\hat{F}(t)]^{\exp\left[\pm 1.96 \sqrt{\text{Var}(\log(-\log(\hat{F}(t))))}\right]}$$

Relating the two (also a NON-PARAMETRIC estimator bc doesn't rely on underlying distr.):

$$\hat{S}(t) = 1 - \hat{F}(t)$$



**Kaplan-Meier Curves**

**Kaplan-Meier estimator:** Combines information from samples with observed failure times and samples that were censored to form a single estimated survival curve meant to reflect the true underlying survival curve S(t).  
• Using right-censored data T\*, we make inference about the distribution of the underlying random variable T.  
• Depends on NON-INFORMATIVE censoring.  
• A **nonparametric estimator** because it does not assume that the data fit any underlying parametric distribution.  
• Similar in spirit to the empirical CDF, but modified to handle censored data  
• Aka **product limit estimator** bc our estimate  $\hat{S}(t)$  is the product of conditional survival probabilities in all intervals that end before time t

$$\hat{S}(t) = \prod_{j: t_j \leq t} \hat{q}_j$$
  
$$\hat{S}(t) = \prod_{j: t_j \leq t} \left(1 - \frac{d_j}{n_j}\right)$$
  
 $\hat{q}_j$  = the conditional probability of surviving past time t<sub>j</sub> given that one has already survived past t<sub>j-1</sub>  
d<sub>j</sub> = failures at time t<sub>j</sub>  
c<sub>j</sub> = censored obs at time t<sub>j</sub>  
n<sub>j</sub> = individuals at risk before time t<sub>j</sub> aka the RISK SET  
**R(t) = {i: T<sub>i</sub> ≥ t}**  
d<sub>j</sub>/n<sub>j</sub> = The proportion of those at risk who fail during the interval  
1- d<sub>j</sub>/n<sub>j</sub> = The proportion of those at risk who survive

Unique failure/censoring time t <sub>j</sub>	Number at risk n <sub>j</sub>	Number of deaths d <sub>j</sub> at t <sub>j</sub>	Number censored c <sub>j</sub> at t <sub>j</sub>	Conditional survival probability $\hat{q}_j$	Kaplan-Meier estimate $\hat{S}(t_{j+1})$
t <sub>0</sub> = 0	n <sub>0</sub> = 21	d <sub>0</sub> = 0	c <sub>0</sub> = 0	$\hat{q}_0 = 1$	$\hat{S}(0) = 1$
t <sub>1</sub> = 2	n <sub>1</sub> = 12	d <sub>1</sub> = 1	c <sub>1</sub> = 0	$\hat{q}_1 = \left(1 - \frac{1}{12}\right)$	$\hat{S}(2) = \hat{q}_1$

**Hazard and cumulative hazard functions**

<b>Survival function</b> $S(t) = \Pr(T > t)$ $S(t) = 1 - F(t)$ $S(t) = e^{-H(t)}$	<b>Hazard function</b> $h(t) = \frac{f(t)}{S(t)}$ $h(t) = \frac{d}{dt} \ln[1 - F(t)]$ $h(t) = -\frac{d}{dt} \ln S(t)$
<b>Cumulative distribution function</b> $F(t) = \Pr(T \leq t)$	<b>Cumulative hazard function</b> $H(t) = \int_0^t h(u) du$ $H(t) = -\ln S(t)$
<b>Probability density function</b> $f(t) = \frac{d}{dt} S(t)$ $f(t) = h(t)S(t)$	

$\frac{d}{dt} F(t) = f(t)$  = Probability density function  
**h(t)** = Hazard function or hazard rate  
$$= \frac{f(t)}{S(t)}$$
  
The instantaneous rate of failure at time t given that you have survived up until time t.  
**H(t)** = Cumulative Hazard Function

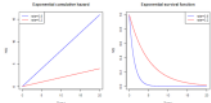
$H(t) = \int_0^t h(u) du$   
The area under the hazard function between the time origin and time t. It sums up the accumulated hazard through time t  
Properties:  
•  $H(0) = 0$ . At time 0, you have accumulated no hazard.  
•  $S(t) = e^{-H(t)}$   
•  $H(t) = -\log S(t)$   
Hazard interval for logit(h(t))  
$$\ln h(t) \pm 1.96 \sqrt{\frac{1}{n} \ln h(t)}$$
  
Log transformed interval for h(t):  
$$h(t) \exp\left[\pm 1.96 \sqrt{\frac{1}{n} \ln h(t)}\right]$$
  
Log-log transformed interval for h(t):  
$$h(t) \exp\left[\pm 1.96 \sqrt{\frac{1}{n} \ln h(t)}\right]$$

**Parametric distributions for time-to-event data**

Distribution	Hazard Function	Cumulative Hazard Function	Survival Function
Exponential	$h(t) = \lambda$	$H(t) = \lambda t$	$S(t) = e^{-\lambda t}$
Weibull	$h(t) = \lambda \gamma (t)^{\gamma-1}$	$H(t) = (\lambda t)^\gamma$	$S(t) = e^{-(\lambda t)^\gamma}$
Log-logistic	$h(t) = \frac{\lambda \gamma (t)^{\gamma-1}}{1 + (\lambda t)^\gamma}$	$H(t) = \log(1 + (\lambda t)^\gamma)$	$S(t) = \frac{1}{1 + (\lambda t)^\gamma}$

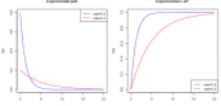
**Exponential Distribution:** Has a constant hazard function

$h(t) = \lambda$   
 $\lambda$  = rate parameter, units of time<sup>-1</sup>  
 $\sigma$  = scale parameter = 1/λ = mean survival time  
$$H(t) = \int_0^t h(u) du = \int_0^t \lambda du = \lambda t$$
  
$$S(t) = e^{-H(t)} = e^{-\lambda t}$$



Note that although the parametric hazard functions are straight lines, the survival functions are curved. A higher rate parameter λ results in worse survival (shorter survival time). Later failure is censored.

pdf	$f(t) = \lambda e^{-\lambda t}$
CDF	$F(t) = 1 - e^{-\lambda t}$
Mean	$E[T] = 1/\lambda$
Variance	$\text{Var}(T) = 1/\lambda^2$
Median	$t_{0.5} = \ln(2)/\lambda$

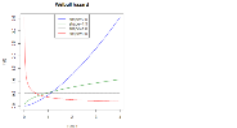


**Weibull Distribution**

The Weibull distribution is a continuous probability distribution.  
 $h(t) = \lambda \gamma (t)^{\gamma-1}$   
 $\lambda$  = rate parameter, units of time<sup>-1</sup>  
 $\gamma$  = shape parameter  
 $\sigma$  = scale parameter = 1/λ  
$$H(t) = \int_0^t \lambda \gamma (u)^{\gamma-1} du = \lambda \int_0^t \gamma u^{\gamma-1} du = (\lambda t)^\gamma$$

$S(t) = \exp(-H(t)) = \exp(-( \lambda t)^\gamma)$   
decreasing if γ < 1  
constant if γ = 1  
increasing if γ > 1  
increasing if γ < 2  
decreasing if γ > 2

When γ = 1, h(t) = λ, and the Weibull distribution simplifies to an exponential distribution.



**Log-log distribution:** The distribution of a random variable whose logarithm has a logistic distribution

$T \sim \text{Log-logistic}(\lambda, \gamma)$   
 $\ln(T) \sim \text{Logistic}(\ln(\lambda), \gamma)$

censoring time $t_j$	$n_j$ during $(t_{j-1}, t_j]$	deaths $d_j$ at $t_j$	$c_j$ at $t_j$	probability $\hat{q}_j$	$P_j, F_j, S_j$
$t_0 = 0$					$t = [0, 2]$ $\hat{S}(t) = 1$
$t_1 = 2$	$t = (0, 2]$ $n_0 = 12$	$d_0 = 1$	$c_1 = 0$	$\hat{q}_0 = \left(1 - \frac{1}{12}\right)$ $\hat{S}(t) = \hat{q}_0$	$t = (2, 3]$ $\hat{S}(t) = \hat{q}_0$
$t_2 = 3$	$t = (2, 3]$ $n_1 = 11$	$d_1 = 0$	$c_2 = 1$	$\hat{q}_1 = 1$ $\hat{S}(t) = \hat{q}_0 \hat{q}_1$	$t = (3, 6]$ $\hat{S}(t) = \hat{q}_0 \hat{q}_1$
$t_3 = 6$	$t = (3, 6]$ $n_2 = 10$	$d_2 = 2$	$c_3 = 0$	$\hat{q}_2 = \left(1 - \frac{2}{10}\right)$ $\hat{S}(t) = \hat{q}_0 \hat{q}_1 \hat{q}_2$	$t = [6, \infty)$ $\hat{S}(t) = \hat{q}_0 \hat{q}_1 \hat{q}_2$

An interpretation of  $S_{\text{hat}}(t)$  is

An estimated 73.3% of hemophiliacs are alive 7 months after primary AIDS diagnosis.

Properties of Kaplan-Meier curve:

- Has a step function appearance.
- Is equal to one up to the first death time.
- It only drops at the time of failure. It does not drop when individuals are censored.
- Drops to 0 if the last event is a death. (Has poor fit in tails.)
- In the absence of censoring, the Kaplan-Meier simplifies to the empirical CDF.
- A Kaplan-Meier plot is the most common figure shown in a paper with time-to-event data.

Greenwood's Formula for CI

$$\widehat{Var} \hat{S}(t) = \left(\hat{S}(t)\right)^2 \sum_{j: t_j \leq t} \frac{d_j}{n_j(n_j - d_j)}$$

$$\hat{S}(t) \pm 1.96 \sqrt{\widehat{Var}(\hat{S}(t))} \quad \longleftrightarrow \quad \text{Standard error is just } \sqrt{\widehat{Var}}$$

Log transformed:

$$\left(\hat{S}(t)\right)^{\log(1.96)} \cdot \left(\hat{S}(t)\right)^{-\log(1.96)}$$

Where  $w = \frac{1}{\sqrt{1.96 \cdot \log(1.96) \cdot n_j(n_j - d_j)}}$

Log-log transformed:

$$\left(\hat{S}(t)\right)^{\log(1.96)} \cdot \left(\hat{S}(t)\right)^{-\log(1.96)}$$

Where

$$w = \frac{1}{\sqrt{\log(\hat{S}(t))}} \sum_{j: t_j \leq t} \frac{d_j}{n_j(n_j - d_j)}$$

## The log-rank test

**Log-rank test:** most commonly used statistical test for comparing the survival functions of two or more independent groups.

- Nonparametric** test whose validity does not depend on any parametric assumptions.
- Not well suited** for the setting where the survival curves cross and the **direction** of the difference changes
- It does not allow us to measure the simultaneous impact of these variables on survival
- It does not borrow information across similar groupings
- We can only stratify on categorical covariates
- We may wish to model the number of positive lymph nodes as a continuous variable
- It does not provide a summary statistic for the effect size

$n_{0j}$  is the number at risk in group 0 at time  $t_j$

$n_{1j}$  is the number at risk in group 1 at time  $t_j$

$d_{0j}$  is the number of failures in group 0 at time  $t_j$

$d_{1j}$  is the number of failures in group 1 at time  $t_j$

Consider each distinct failure time  $t_j$

	At risk just before $t_j$	Fail at $t_j$	Survived past $t_j$
Group 0	$n_{0j}$	$d_{0j}$	$n_{0j} - d_{0j}$
Group 1	$n_{1j}$	$d_{1j}$	$n_{1j} - d_{1j}$
Total	$n_j = n_{0j} + n_{1j}$	$d_j = d_{0j} + d_{1j}$	$n_j - d_j$

**Null hypothesis  $H_0$**

- $S_0(t) = S_1(t)$  for all  $t$
- $S_0(t) > S_1(t)$
- The survival probabilities are equal at all times  $t$

**Alternative hypothesis  $H_A$  or  $H_1$**

- $S_0(t) \neq S_1(t)$  for some  $t$
- $S_0(t) < S_1(t)$
- The survival probabilities are different at some time(s)  $t$

Under  $H_0$ , the expected number of failures in group 0 at time  $t_j$  is:

$$E_j = \frac{n_{0j}}{n_j} d_j$$

This is the proportion of people at risk who are in group 0, multiplied by the total number of events at that time.

$O_j = d_{0j}$

We standardize the difference between  $O_j$  and  $E_j$  by the variance  $V_j$  under  $H_0$

$$V_j = \frac{n_{0j} n_{1j} (d_j - E_j)}{n_j^2 (n_j - 1)}$$

$$n = \sum_{j=1}^J n_j$$

$$Z = \sum_{j=1}^J \frac{O_j - E_j}{\sqrt{V_j}}$$

$$V = \sum_{j=1}^J V_j$$

The log-rank test statistic  $Z$  is related to:

$$Z = \frac{O - E}{\sqrt{V}}$$

**Weighted log-rank test**

Let  $w_1 \geq 0, w_2 \geq 0, \dots, w_J \geq 0$  be known constants (weights). Then the weighted log-rank test is given by:

$$Z_w = \frac{\sum_{j=1}^J w_j (O_j - E_j)}{\sqrt{\sum_{j=1}^J w_j^2 V_j}}$$

**Generalized Wilcoxon Test/Gehan-Breslow Test:** Choosing  $W_j = n_j$

When group difference is large early on but then decreases over time: Generalized Wilcoxon test statistic > standard log-rank test statistic  
When group difference is small early on but then increases over time: Generalized Wilcoxon test statistic < standard log-rank test statistic

**Stratified log-rank test:**

The null hypothesis is that the survival functions of the two groups of interest (e.g. treated and untreated) are *the same within each stratum*.

Does not assume that the survival functions in different strata are the same, which is important because there might be differences in the shape of the survival function across strata. Instead we are estimating the treatment differences within each stratum and then pooling these differences across strata.

Pros:

- Including a categorical covariate allows us to "explain" some of the variability observed between individuals.
- Looking within groups allows us to more precisely isolate the effect of treatment.
- There may be confounding. By looking within groups, we address these imbalances

$$H_{0S}: S_0^{(k)}(t) = S_1^{(k)}(t)$$

$$O_k = \sum_{j=1}^J O_{kj}$$

$$E_k = \sum_{j=1}^J E_{kj}$$

$$V_k = \sum_{j=1}^J V_{kj}$$

## Benefits of stratification

- Survival may vary a lot across levels of the stratifying factors

- For non-randomized studies, the two groups may be imbalanced with respect to these factors, leading to confounding
- Stratification can address this confounding

- For randomized trials, the two groups are expected to be balanced so confounding is not a concern and stratification is not necessary for validity
- But adjusting for other covariates that are predictive of survival may decrease variability in the population and improve our ability to detect differences across groups

## Statistical inference



**Log-log distribution:** The distribution of a random variable whose logarithm has a logistic distribution

$$T \sim \text{Log-logistic}(\lambda, \gamma)$$

$$\log(T) \sim \text{Logistic}(\lambda, \gamma)$$

$$h(t) = \frac{\lambda \gamma (t)^{\gamma-1}}{1 + (t^\gamma)^\gamma}$$

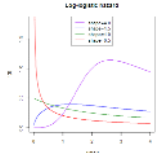
$$\lambda = \text{rate parameter, units of time}^{-1}$$

$$\gamma = \text{shape parameter} = 1/\lambda$$

$$\sigma = \text{scale parameter} = 1/\lambda$$

$$h(t) \text{ is } \begin{cases} \text{decreasing from } \sigma_0 & \text{if } \gamma < 1 \\ \text{decreasing from } \lambda & \text{if } \gamma = 1 \\ \text{increasing then decreasing} & \text{if } \gamma > 1 \end{cases}$$

$$\begin{aligned} H(t) &= -\log\left(\frac{1}{H(t)}\right) \\ S(t) &= \exp\left(-\log\left(\frac{1}{H(t)}\right)\right) \\ &= \frac{1}{1 + (t^\gamma)^\gamma} \end{aligned}$$



## Maximum likelihood

A general method for using data to estimate parameters such as those used in the exponential, Weibull, and log-logistic distributions.

- If  $T_j^*$  is a failure time ( $\delta_j = 1$ ), the probability of observing this failure time can be expressed by the pdf at time  $T_j^*$ .
- If  $T_j^*$  is a censoring time ( $\delta_j = 0$ ), then we do not directly observe when this person failed, but we do observe that this person survived at least up until  $T_j^*$ . Thus, the probability of observing this censoring time can be expressed by the survival function at time  $T_j^*$ .

Consider a distribution with a single parameter  $\lambda$ . The **likelihood function**  $L(\lambda)$  evaluated at  $\lambda$  can be conveniently expressed as follows:

$$L(\lambda) = \prod_{i=1}^n [\lambda (T_i^*)^{\lambda-1}]^{\delta_i} S(T_i^* | \lambda)$$

e.g. likelihood function for exponential distribution

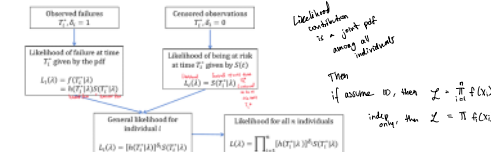
$$L(\lambda) = \prod_{i=1}^n \lambda e^{-\lambda T_i^*}$$

$$\hat{\lambda} = \text{MLE}$$

e.g. MLE for exponential distribution, also incidence rate

$$\hat{\lambda} = \frac{\sum_{i=1}^n \delta_i}{\sum_{i=1}^n T_i^*}$$

## Likelihood for right-censored data



$$\text{General Likelihood: } L(\lambda) = \prod_{i=1}^n [\lambda (T_i^*)^{\lambda-1}]^{\delta_i} S(T_i^* | \lambda)$$

Distribution	Hazard Function	Survival Function	Likelihood Function
Exponential	$h(t) = \lambda$	$S(t) = e^{-\lambda t}$	$L(\lambda) = \prod_{i=1}^n \lambda e^{-\lambda T_i^*}$

$$\text{Weibull} \quad h(t) = \lambda \gamma (t)^{\gamma-1} \quad S(t) = e^{-(\lambda t)^\gamma} \quad L(\lambda, \gamma) = \prod_{i=1}^n [\lambda \gamma (T_i^*)^{\gamma-1}]^{\delta_i} e^{-(\lambda T_i^*)^\gamma}$$

$$\text{Log-logistic} \quad h(t) = \frac{\lambda \gamma (t)^{\gamma-1}}{1 + (t^\gamma)^\gamma} \quad S(t) = \frac{1}{1 + (t^\gamma)^\gamma} \quad L(\lambda, \gamma) = \prod_{i=1}^n \left[ \frac{\lambda \gamma (T_i^*)^{\gamma-1}}{1 + (T_i^*)^\gamma} \right]^{\delta_i} \left[ \frac{1}{1 + (T_i^*)^\gamma} \right]^{1-\delta_i}$$

**Nelson-Aalen / Breslow / Fleming-Harrington Estimator:** Another non-parametric estimator of the survival function. Sort of like Riemann sums

$$\hat{h}(t_j) = \frac{d_j}{n_j(t_j - t_{j-1})}$$

$$\hat{H}(t) = \sum_{t_j \leq t} \hat{h}(t_j)$$

$$\widehat{Var}(\hat{H}(t)) = \sum_{t_j \leq t} \frac{d_j}{n_j^2}$$

Breslow estimator:

$$\hat{S}(t) = \exp(-\hat{H}(t))$$

A log-transformed 95% confidence interval for  $H(t)$  is:

$$\hat{H}(t) \exp\left(\pm 1.96 \sqrt{\widehat{Var}(\hat{H}(t))}\right)$$

$$(e^{-\hat{H}(t)}, e^{-\hat{H}(t)})$$

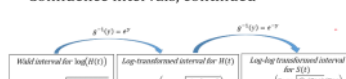
## Point estimation



## Confidence intervals



## Confidence intervals, continued



The diagram illustrates the process of hypothesis testing, showing the relationship between a population and a sample.

**Population (Left):** We want to know if these (represented by 10 stick figures). The population is labeled  $S(T)$  (True control function). The goal is to "Test whether  $\hat{g}(T) = g_0(T)$ ".

**Sampling Process:** Random selection (indicated by a grey arrow) leads to a sample.

**Sample (Right):** We have those to work with (represented by 3 stick figures). The sample is labeled  $\hat{S}(T)$  (Estimated control function). The goal is to "Estimate  $\hat{g}_1(T)$  and  $\hat{g}_0(T)$  for two independent populations" and "Statistic".

**Statistical Inference:** A grey arrow labeled "Hypothesis testing Inference" points from the sample back to the population.

**Additional Information:** A note on the right states: "Sample of right-censored survival times  $(T^*, \delta)$  from two independent populations".

The diagram illustrates the complementary log-log (clog-log) transformation. It consists of three boxes arranged horizontally, connected by curved arrows. The top row of boxes contains the following text:

- Wild interval for  $\log(H(t))$
- Log-transformed interval for  $H(t)$
- Log-log transformed interval for  $S(t)$

The bottom row of boxes contains the following text:

- $\log(H(t)) \pm 1.96 \sqrt{\text{Var}(\log(H(t)))}$
- $H(t) \exp\left(\frac{\pm 1.96 \sqrt{\text{Var}(H(t))}}{H(t)}\right)$
- $S(t) \exp\left(\pm 1.96 \sqrt{\text{Var}(S(t))} \log(H(t))\right)$

Arrows indicate the transformations between these intervals:

- A blue arrow from the top-left box to the top-middle box is labeled  $g^{-1}(x) = e^x$ .
- A blue arrow from the top-middle box to the top-right box is labeled  $g^{-1}(x) = e^{-x}$ .
- A blue arrow from the bottom-left box to the bottom-middle box is labeled  $g(x) = \log(x)$ .
- A blue arrow from the bottom-middle box to the bottom-right box is labeled  $g(x) = -\log(x)$ .

At the bottom, the text "Complementary log-log:  $\log(-\log(x))$ " is displayed.