✓ <u>Survival Analysis</u>

Survival Analysis is a collection of statistical procedures for ~~better~~ data analysis for which, the outcome variable of interest is time untill an event occures, By time we mean years, months, weeks on days from the begining of follow up of an individual untill and event occures, alternatively time can reffer to the age of an individual, when and event occuns. By event we mean death disease x incidence, relapse from remission, recovery on, any designated experience of interest that may happen to an individual.

<u>Censoring</u>

Most survival analysis must consider a key analytical problem called censoring. Censoring occures when we have some information about individual survival time, but, we do not know the survival time exactly.

<u>Example:-</u> As a simple example consider leukemia patient followed until they go out of remission. If for a given patient, the study ends while the patient still in remission. (ie. doesn't get the event) that patient survival time is consider censored. For this person, time is atleast as long as the period that the person has been followed, but if the person goes out of remission after the study ends, we don't know the comple survival time.

(doesn't get the event → remission)

Q. There are generally three reasons for censoring. They are —

① A person doesn't experience the event before the study ends.

② A person is lost to follow up during the study period.

③ A person withdraws from the study.

* Types of Cencering :-

There are three types of cencering :

① Right censored :

True survival time is equal to on greaten then observed survival rate time.

② Left censored :

True survival time is less then -on equal to the observed survival time.

③ Interval censored :

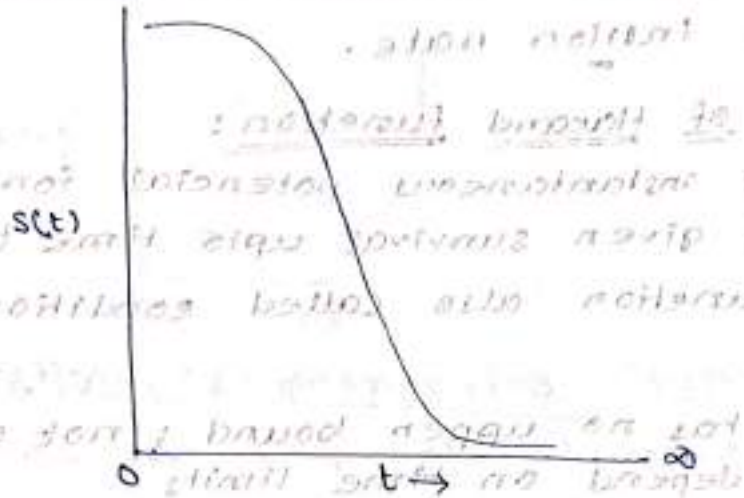True survival time is with in a know time survival.

<u>Surviver function</u> :

The survivar function S(t) this the probability that a person survive longer then some specifiee time t i.e. S(t) gives the probability that the random variable T exists the specified time t.

Theorytically as t ranges from 0 to infinite the survivor function can be ratte as a smooth curve As clastrated by the group where t identyfies the X exis all survival function have the following charectenistic :—



S(t)

① They are nonincreasing, i.e. they head down word as t increases.

② At time $t=0$, $S(t) = S(0) = 1$, i.e. at the start of the study since know one can gotten have even the probability of surviving past time 0 is one

③ At time $t = \infty$, $S(t) = S(\infty) = 0$, i.e. theorytically if the study period inencrease without limit eventually nobody would survivoun. so the survival past fall to 0.

● <u>Hazand function</u> :—

The Hazand function denoted by

$$h(t) = \int \lim_{t \to 0} \Delta t$$

The hazand function h(t) gives the instantaneous potensial per unit time for the even to occun, given that the individual has survive upto time t. In contras to the survival function which

foceuse on not failing on hazard: function focus [9]
on falling i.e. on the event occuring. Thus in
some sense, can be considered as given the opposite
the hazard function
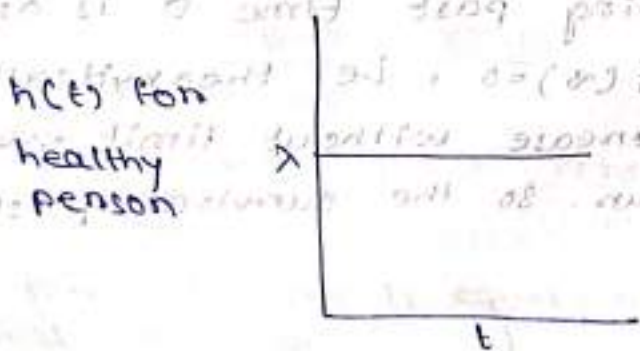site
time of the information, survival function.
giving by the

In the hazard formula, the conditional
probability gives the probability that a person
survival time T WE line, interval between T ant+ot,
the time
given that survival time is greater then on
equal to t. Because of the 'given' sign means
hence the hazard function is sometimes
conditional failion rate.
or

• **Property of Hazard function:**

① H(t) gives instantaneous potencial for even to
occur to given survival upto time t.

② Hazand function also called conditional failon
rate.

③ H(t) ⩾ 0; has no upper bound; not a prob;
depend on time limits

## Types of Hazard function:

① Hazard function is constant: The following graph
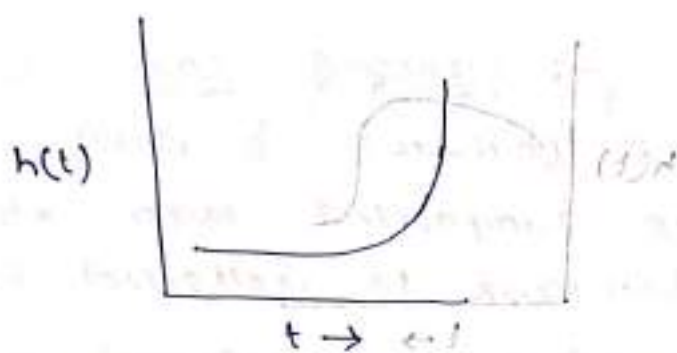given shows constant hazard for study a
healthy person.



h(t) for
healthy
person

In this graph no matter what value of t
specified h(t) - the some value. Hence h(t) + λ.
when the hazard function is constant we
say the survival model is exponential.

② Hazard function in increasing overtime:
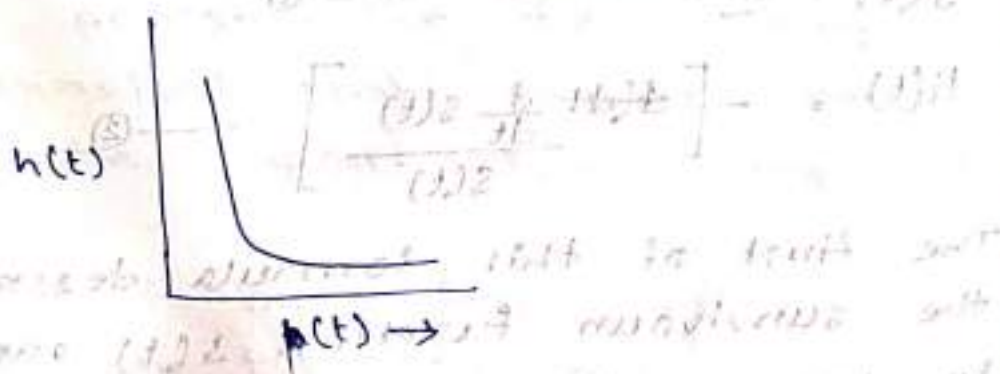The following graph shows a hazard function

i.e. increasing overtime. And example of this kind of graph is called and increasing weibull model. Such a graph might be expected for leukemia patients not responding to treatment where the event of interest is death. As survival time increases for such a patient and as the prognosis accordingly worsens, the patients potential for dying of the disease also increases.



h(t)
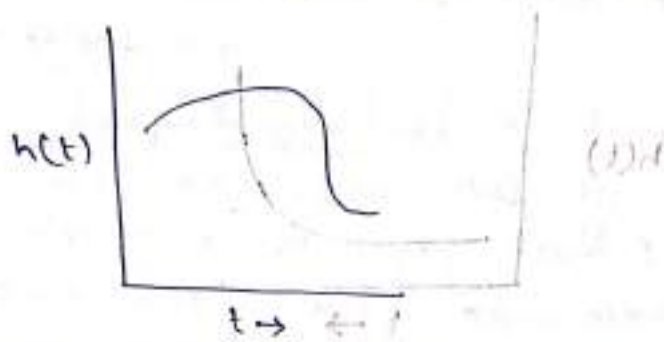
t →

③ Hazard function is decreasing overtime:

In the following graph the hazard function is decreasing overtime. An example of this kind of graph is called (a) decreasing weibull. Such a graph might be expected when the event is death in persons who are recovering from surgery, because the potential for dying after surgery usually decreases as the time after surgery increases.



h(t)

t →

④ **Hazard Function i.e. first increasing and then decreasing:**

The following graph shows a hazard function i.e. first increasing and then decreasing. An example of this type of graph is the log normal survival model. We can expect such a graph for Tuberculosis patients, since there potential for dyeing increases early in the disease and decreases later.



**Relationship of ~~S(t)~~ S(t) and H(t):-**

There is a cleanly defined relationship between the survival function (s(t)) and Hazard Function (H(t). In fact if one knows the form of S(t), one can deplve the corresponding H(t) and the vice versa. The relationship between S(t) and H(t) can be expressed equivalently in either of two calculus formula shown in following:

$$S(t) = e^{-\int_0^t h(u)d(u)} \quad —— ①$$

$$H(t) = -\left[ \frac{\frac{d}{dt}s(t)}{s(t)} \right] \quad —— ②$$

The first of this formula describes how the survivour function S(t) can be written in terms of an integral involving the Hazard Function. The formula

says that $s(t)$ equals the exponential of the -ve integral of the Hazard Function between integration units of 0 and t. The

The second formula describes how the Hazard Function $H(t)$ can be written in terms of derivative involving the survivoun function. These formula says that $H(t)$ equals — the derivative of $s(t)$ with respected to t devided by $s(t)$.

**Goals of Survival Analysis :-**
The basic goals of survival analysis are -
① To Estimate and interpret survivoun and / or hazard function of survival data.

② To compane survivoun and / on hazand function

③ To esses the relationship of explanatory variab to survival time.

**Kaplan-Meier Survival Curves :-**
To estimate the survival probability at a given time we make use of the risk set at that time to include the information we have on a censored person upto the time censorship, rather than simply through a away all the information on a censored person. The extual computation of such survival prob. can be carried out using the Kaplan-meien method (KM method)

The general formula for a KM survival prob. at failure time $t_{(f)}$ eat is given by

$$\hat{S}(t_{(f)}) = \hat{S}(t_{(f-1)}) \times \hat{P}_n(T > t_{(f)}) / (T \geq t_{(f)})$$

This formula gives the prob. of surviving past the previous failure time $t_{(f-1)}$, multiplied by the conditional probability of surviving past time $t_{(f)}$, given survival to atleast time $t_{(f)}$. The above KM formula can also be expressed as a product limit if we substitute for the survival probability as $\hat{S}(t_{(f-1)})$, the product of all fraction that estimate the conditional prob. for failure times $t_{(f-1)}$ and earlier.

i.e. $\hat{S}(t_{(f-1)}) = \prod_{i=1}^{f-1} [\hat{P}_n(\hat{T} > t_{(i)} / T \geq t_{(i)})]$

① The remission times (weeks) for two groups of lucomia patients :

Group 1 (N = 21)        Group 2 (N = 21)

Treatment               Placebo

6, 6, 6, 6, 7, 10, 13,      1, 1, 2, 2, 3, 4, 4,
16, 22, 23, 6+, 9+,        5, 5, 8, 8, 8, 8,
10+, 11+, 17+, 19+,       11, 11, 12, 12, 15,
20+, 25+, 32+, 32+,       17, 22, 23
34+, 35+

| $t_{(f)}$ | $n_f$ | $m_f$ | $q_f$ | $\hat{S}(t_f)$ |
|---|---|---|---|---|
| 0 | 21 | 0 | 0 | |
| 6 | 21 | 3 | 1 | 1 × 18/21 = 0.8571 |
| 7 | 17 | 1 | 1 | 0.8571 × 16/17 = 0.8067 |
| 10 | 15 | 1 | 2 | 0.8067 × 14/15 = 0.7529 |
| 13 | 12 | 1 | 6 | 0.7529 × 11/12 = 0.6901 |
| 16 | 11 | 1 | 3 | 0.6901 × 10/11 = 0.6274 |
| 22 | 7 | 1 | 0 | 0.6274 × 6/7 = 0.5371 |
| 23 | 6 | 1 | 5 | 0.5371 × 5/6 = 0.4482 |
| >23 | | | | |

| $t_{(f)}$ | $n_f$ | $m_f$ | $q_f$ | $\hat{S}(t_f)$ |
|---|---|---|---|---|
| 0 | 21 | 0 | 0 | $\rightarrow$ 1 |
| 1 | 21 | 2 | 0 | $\rightarrow$ $1 \times 19/21 =$ |
| 2 | 19 | 2 | 0 | $\rightarrow$ $19/21 \times 17/19 = \dfrac{17}{21}$ |
| 3 | 17 | 1 | 0 | $\rightarrow$ $17/21 \times 16/17 = \dfrac{16}{21}$ |
| 4 | 16 | 2 | 0 | $\rightarrow$ $16/21 \times 14/16 = \dfrac{14}{21}$ |
| 5 | 14 | 2 | 0 | $\rightarrow$ $14/21 \times 12/14 = \dfrac{12}{21}$ |
| 8 | 12 | 4 | 0 | $\rightarrow$ $12/21 \times 8/12 = \dfrac{8}{21}$ |
| 11 | 8 | 2 | 0 | $\rightarrow$ $8/21 \times 6/8 = 6/21$ |
| 12 | 6 | 2 | 0 | $\rightarrow$ $6/21 \times 4/6 = 4/21$ |
| 15 | 4 | 1 | 0 | $\rightarrow$ $4/21 \times 3/4 = 3/21$ |
| 17 | 3 | 1 | 0 | $\rightarrow$ $3/21 \times 2/3 = 2/21$ |
| 22 | 2 | 1 | 0 | $\rightarrow$ $2/21 \times 1/2 = 1/21$ |
| 23 | 1 | 1 | 0 | $\rightarrow$ ~~~~~ 0 |

# Unit - 4

## Log - Rank Test :-

(The Log-Rank test is a large sample $\chi^2$-test that uses as its test criterian a statistic that provides and overall comparison of the k-m curbs being compared) This statistic, like many other statistic used in other kinds of $\chi^2$-test, makes use of observed V/s expected, cell counts over chatagories of outcomes. The (catagories for the Log-Rank statistics are defined by each of the ordered failure times for the entire set of data being analysed (Here, for each order failure time $t_{(f)}$, in the entire set of data, we need the nos of subjects ($m_{if}$) failing at that time seperately by group (i), followed by the no. of subjects ($n_{if}$) in the risk set at that time, also, seperately by group. Now, expected cell counts for each group is calculated by each

| $t_{(i)}$ | $m_{1f}$ | $m_{2f}$ | $n_{1f}$ | $n_{2f}$ | $q_{1f}$ | $q_{2f}$ |
|---|---|---|---|---|---|---|
| 1 | 0 | 2 | 121 | 21 | 0 | 0 |
| 2 | 0 | 2 | 21 | 19 | 0 | 0 |
| 3 | 0 | | 21 | 17 | 0 | 0 |
| 4 | 0 | 2 | 21 | 16 | 0 | 0 |
| 5 | 0 | 2 | 21 | 16 | 0 | 0 |
| 6 | 3 | 0 | 21 | 14 | 0 | 0 |
| 7 | 1 | | 17 | 12 | 1 | 0 |
| 8 | 0 | 4 | 16 | 12 | 1 | 0 |
| 10 | 1 | 0 | 15 | 12 | 2 | 0 |
| 11 | 0 | 2 | 13 | 8 | 2 | 0 |
| 12 | 0 | 2 | 12 | 8 | 0 | 0 |
| 13 | 1 | 0 | 12 | 6 | 1 | 0 |
| 15 | 0 | 1 | 11 | 4 | 0 | 0 |
| 16 | 1 | 0 | 11 | 4 | 0 | 0 |
| 17 | 0 | 1 | 10 | 3 | 3 | 0 |
| 22 | 1 | 1 | 7 | 2 | 0 | 0 |
| 23 | 1 | | | | 0 | |

Now, expected cell counts for each groups is calculated by following formula:-

$$e_{1f} = \left(\frac{n_{1f}}{n_{1f} + n_{2f}}\right) \times (m_{1f} + m_{2f})$$

$$e_{2f} = \left(\frac{n_{2f}}{n_{1f} + n_{2f}}\right) \times (m_{1f} + m_{2f})$$

when, two groups are being compared the Log-Rank test statistics is form using the sum of the observed - expected counts over all failure times for one of the two groups. For the two groups case the log rank test statistic is defined by

Log-rank test statistic,

$$= \frac{(O_i - E_i)^2}{v(O_i - E_i)} \, , \quad i = 1 \ \text{or} \ 2$$

where.

$$v(O_i - E_i)$$

$$= \sum_{f} \frac{n_{1f}\, n_{2f}\, (m_{1f} + m_{2f})\, (n_{1f} + n_{2f} - m_{1f} - m_{2f})}{(m_{if} + n_{if})^2\, (n_{1f} + n_{2f} - 1)} \quad ; \quad i = 1 \text{ on } 2$$

The null hypothesis being tested is that there is no overwall difference between the 2 survival curves. Under these null hypothesis the log-rank test statistics is approximately $\chi^2$ with one degrees of freedom. i.e.

$$\chi^2 \simeq \sum_{i=1}^{2} \frac{(O_i - E_i)^2}{E_i} = \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2}$$

Now, if the calculated $\chi^2 > \chi^2_{0.05}$ at $\alpha$ that we reject the null hypothesis otherwise we accept it.

Date = 15/07

Formula for the Log-Rank Statistics for several Groups :-

The Log-Rank test can also be used to compare three on more survival curves. The null hypothesis for these more general situation is that all survival curves are the same. Here we need to calculate the variances and covariances of $O_i - E_i$. For $i = 1, 2, \ldots, G$ and $f = 1, 2, \ldots, k$ where, $G = $ no. of groups and $K = $ no. of failure times, $N_{if} = $ no. at risk in the $i$th group and the $f$th and the $N_{if} = $ observed no. of failures of $i$th group and $f$th orden to failure time and in $i$th group at $f$th ordened failure time $= \frac{N_{if}}{N_{1f} + N_{if}}$

$$\Rightarrow n_f = \sum_{i=1}^{q} n_{if}$$

$$m_f = \sum_{i=1}^{q}$$

$$V(O_i - E_i) = \sum_{f=1}^{k} \left( \frac{n_{if}(n_f - n_{if}) m_{if}(n_f - m_f)}{n_f^2 (n_f - 1)} \right)$$

$$Cov(O_i - E_i, O_L - E_L) = \sum_{f=1}^{k} \left( \frac{-n_{if}(n_{iL} \cdot m_f (n_f - m_f)}{n_f^2 (n_f - 1)} \right)$$

$$d = \left( O_1 - E_1, O_2 - E_2, \ldots, O_{G-1} - E_{G-1} \right)^T$$

$$V = \left( \begin{matrix} (V_{iL}) \\ \vdots \end{matrix} \right) \quad \text{where, } V_{ii} = V(O_i - E_i) \text{ and}$$

$$V_{iL} = Cov(O_i - E_i, O_L - E_L)$$

$$i = 1, 2, \ldots, G-1$$
$$L = 1, 2, \ldots, G-1$$

Then, the log-Rank Statistic is given by the matrix product formula :

Log-Rank Statistic $= d' V^{-1} d$ which has approximately a $\chi^2$-distribution with $G-1$ degrees of freedom under the null hypothesis that all $G$ groups have common survival curve.

<u>Alternatives to the Log-Rank Test</u> :-

There are several alternatives to the Log-Rank Test. The Log-Rank uses $O_i - E_i = \sum m_{if} - e_{if}$
$i = $ no. of groups
$f = f^{th}$ failure time

This simple some gives the same weight namely unity to each failure time when combining observed - expected failure in each group. ~~waiting the~~

Here the test statistic is ,

$$\frac{\left( \sum w(t_{(f)}) (m_{if} - e_{if}) \right)^2}{V \left( \sum w(t_{(f)}) (m_{if} - e_{if}) \right)}$$

where, $i = 1, 2$
$f = f^{th}$ failure time
$w(t_{(f)}) = $ weight at $f^{th}$ failure time

The Wilcoxon, Tarone-Ware, Paro, Flamington-Harington. Test statistics are variation of the Log-Rank test statistics are derived by applying different weights at the $f^{th}$ failure time. The Wilcoxon Test weights the observed - expected score at time $t_{(f)}$ by the no. at risk, $nf$, overall groups at time $t_{(f)}$. Thus the Wilcoxon-Test places more emphasise on the information at the begining of the survival curve at the no. of risk at large allowing early failure to receive more weight then leter fellow. This type of weighting may be used to esses wheathen the effect of a theatment on survival strongest in earlier phases of administration and tens to be less effective over time.

The Tarone-Ware test statistics also applies more way to the early failure times by weighting the observed-expected score at time $t_{(f)}$ by the square root of the no, at risk $\sqrt{nf}$.

The Paro Test weightes the effect failure time the survival estimatise $S(t_{(f)})$ calculated overall groups combine

The Flamington-Harington Test uses Kaplan-Myeien survival estimate $\hat{S}(t)$ overall groups to calculate its weights for the $f^{th}$ failure time,

$$\hat{S}(t_{(f-1)})^p \left[1 - \hat{S}(t_{(f-1)})\right]^q$$

The Flamington-Harington Test allowes the most flexibility in terms of the coice of ways beta weights because the user provides the values of $p$ and $q$. For example, $p=1$ and $q=0$ and $W_p = \hat{S} t_{(f-1)}$

which gives more weights for the earlier survival times when $\hat{S}(t_{(i-1)})$ is closed to 1. However if $p=0$ and $q=1$ then $w_i = 1 - \hat{S} t_{(i-1)}$ in which case the letus survival time receives more weight. If $p=0$ and $q=0$ then $w_i = 1$ and the Flamington-Harington test reduces to Log-Rank Test.

## Cox-Proportional Hazard Model

The Cox-Proportional Hazard model is defined by $h(t, X) = h_0(t) \cdot e^{\sum_i \beta_i X_i}$. These model gives and expression for the hazard at time t for an individual with a given specification of a set of exploratory variables denoted by the $X$. i.e. the $X$ represents a collection of predicted variable i.e. being model too predict an individuals hazard. The Cox-Model formula says that the hazard a time 't' is the product of two quantities. The 1st of these $h_0(t)$ is called the baseline hazard function. The 2nd quantity is the exponential expression to the linear sum of $\beta_i X_i$, where, the sum is over the p explanatory X variables. And important feature of these formula, which concern the proportional hazards assumption is that the baseline hazard is a function of t but doesn't involve the X's. In contrast the exponential expression below here involves X but doesn't involve t.

The X are called time dependent X's. The Cox-model formula has the property that if all the X's are equal to zero. The formula reduces to the baseline hazard function. i.e. the exponential part of the formula become $e^0$, which is 1. These property of the Cox-model is the reason $h_0(t)$ is called the baseline function or from a slightly different puspective the Cox-model reduces to the baseline hazard when no X's are in the model. Thus $h_0(t)$ may be considend as a starting or baseline version of a hazard function priod to considening any of the exces. Another impontand property of the Cox-model is that the baseline hazard @ $h_0(t)$ is an unspecified function. It is this property that makes the Cox-model is semi-parametric model.

Why Cox-Proportional Hazard Model is popular?

A key neason for the popularity of the Cox-model is that even though the baselin hazard is not specified, neasonable good estimate, negression coefficient and hazard naHos of interest and adjusted survival cunve can be obtain for a wide varialty of data situations. Another way of saying this is that the Cox-PH model is a "Robust model". so that the nesults from using the Cox model will closely apphroximately the

results for the correct parametric model, we would prefer to use a parametric model if we were sure of the correct model. Thus when in doubt the cox-model will give reliable enough results so that it is a "safe" choice of model, and the user doesn't need to worry about wheather the wrong parametric model is choosen.

**maximum likelihood estimation of the cox-proportional hazard model :-**

The maximum likelihood estimators of the cox-model parameters are derived by maximimising a likelihood function, usually denoted by $L$. The maximum likelihood function is a mathematical expresion which describes the joint prob. of obtaining the data actually observed on the subjects in the study as a function of unknown parameters ( the $\beta$'s) in the model being considered. $L$ is sometimes written no notationally as $L(\beta)$ where $\beta$ denots the collection of unknown parameters.

The formula for the cox-model likelihood function is actually called a partial function rather than a complete likelihood function. The term partial likelihood is used because the likelihood formula considers prob. only for those subjects who fail and doesn't explicitly considered prob. for those subjects who are censored. thus the likelihood for the cox-model doesn't consider prob. for all subjects and so it is called a partial

likelihood. In particular the partial likelihood can be written as the product of several one for each of say k failure times. Thus at the $f^{th}$ failure time $L_f$ denotes the likelihood of failing at this time given survival upto this time. Thus $L = L_1 \times L_2 \times L_3 \times \cdots \times L_k$

$$= \prod_{j=1}^{k} L_j$$

Thus although the partial likelihood focuses on subjects who failed the survival time information prior to censorship to use those subject who are censored. i.e. a person who is censored after the $f^{th}$ failure time is part of the risk set used to compute $L_f$ even though this person is censored later. Once the likelihood function for a given model the maximization process is carried out by taking partial derivatives of $\log(L)$ with respect to each parameter in the model and then solving a system of equn.

$$\frac{\partial \log(L)}{\partial \beta_i} = 0 \quad , \quad i = 1, 2, \ldots, p$$

This solution is carried out using iteration. i.e. the solution is based on a stepwise manner which starts a guessed value for the solution and then successively modifieds the guessed value until a solution is finally obtain.

Hazard Ratio :-

In general a hazard ratio is define as the hazard for one individual devided by the hazard for a different individual. The two individual being compared can be distinguis by the values of the set of predictorr, i.e,

the x's. We can write the hazard ratio
as the estimate of $\underline{h(t, x^*)}$ divided by
the estimate of $h(t, x)$. where $x^*$ denotes
the set of predictors for one individual
and $x$ denotes the set of predictors
for other individuals. Thus, $\hat{HR} = \dfrac{\hat{h}(t, x^*)}{\hat{h}(t, x)}$

<u>The meaning of the PH assumption</u> :- ✓

The PH assumption requires that the HR is
constant overtime, or equivalently that the
hazard for one individual is proportional
to the hazard for any other individual,
where the proportionality (ant is independent
of time. We know that To understand the
PH assumption we consider the formula
for the HR that compares two different
specification $x^*$ and $x$ for the explanatory
variable use in the cox - model. Thus for

$$\hat{HR} = \frac{\hat{h}(t, x^*)}{\hat{h}(t, x)}$$

$$= \frac{\hat{h}_0(t) \times e^{\sum \hat{\beta}_i \cdot x_i^*}}{h_0(t) \times e^{\sum \hat{\beta}_i \cdot x_i}}$$

$$= e^{\sum \beta_i (x_i^* - x_i)} \quad \text{doesn't involve } t$$

$$= \hat{\Theta} \text{ (say)}$$

Thus, once the model is fitted and the
value of $x^*$ and $x$ are specified the
value of the exponential expression
for the estimated hazard ratio is a cont.
which doesn't depend time. It we denote
this cont by $\hat{\Theta}$ then we can write the
hazard ratio as shown above. This is a

a mathamatical expression which states the proportional Hazand assumption. (99)