# statstutor

# community project

encouraging academics to share statistics support resources

stcp-marshall-survival

## Survival Analysis

Survival analysis is concerned with data where we measure the time to some event and the outcome of interest is the time to an event. Commonly the event is death (hence the name survival analysis), but it can be other outcomes.

**DATA:** Worcester Heart Attack Study data from Dr. Robert J. Goldberg of the Department of Cardiology at the University of Massachusetts Medical School

**DESCRIPTIVE ABSTRACT:** The main goal of this study is to describe factors associated with trends over time in the incidence and survival rates following hospital admission for acute myocardial infarction (MI). Data have been collected during thirteen 1-year periods beginning in 1975 and extending through 2001 on all MI patients admitted to hospitals in the Worcester, Massachusetts Standard Metropolitan Statistical Area.

## Producing a Kaplan-Meier Plot

A Kaplan-Meier plot displays survivals curves (cumulative probability of an individual remaining alive/ disease free etc. during a unit of time). Note: The cumulative survival probability is the product of the survival probabilities up to that point in time.

ANALYSE → SURVIVAL → KAPLAN-MEIER and select the following options:



Define event: Tell SPSS what number defines an event. Death is indicated by 1 for this study.

If you are interested in the difference between two groups enter the grouping variable into the factor box: here we are interested in the difference between males and females so enter gender into the factor box. If not then leave the box blank.

©Ellen Marshall
University of Sheffield

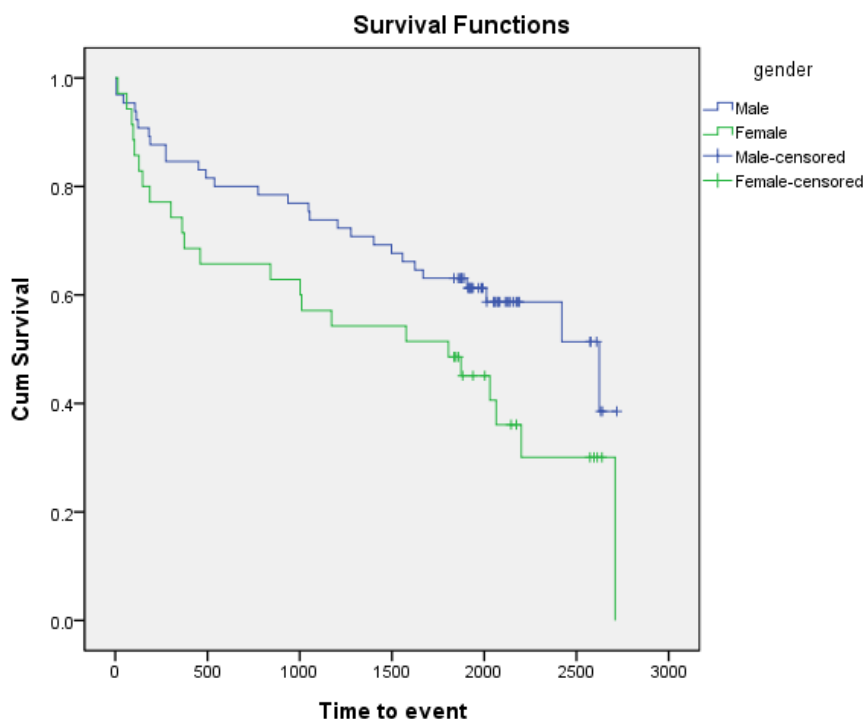Reviewer: Chris Knox
University of Sheffield

Summary statistics for the two groups: survival times should be summarised using the median time to event (shown in the following figure). Estimated time until death is 2624 days for males and 1806 days for females following admission for acute myocardial infarction.

**Means and Medians for Survival Time**

| gender | Mean[a] | | | | Median | | | |
|---|---|---|---|---|---|---|---|---|
| | Estimate | Std. Error | 95% Confidence Interval | | Estimate | Std. Error | 95% Confidence Interval | |
| | | | Lower Bound | Upper Bound | | | Lower Bound | Upper Bound |
| Male | 1907.423 | 126.011 | 1660.442 | 2154.405 | 2624.000 | 392.487 | 1854.726 | 3393.274 |
| Female | 1475.214 | 185.790 | 1111.066 | 1839.363 | 1806.000 | 520.636 | 785.554 | 2826.446 |
| Overall | 1750.270 | 105.951 | 1542.607 | 1957.934 | 2201.000 | 251.678 | 1707.712 | 2694.288 |

a. Estimation is limited to the largest survival time if it is censored.

The Kaplan-Meier plot shows that the survival probability is lower for females at all time points so they are less likely to survive.  Censoring means that an individual is still alive at the end of the study or that they withdrew from the study at that point in time.

**The Log Rank Test**

The log-rank test investigates the hypothesis that there is no difference in survival times between the groups studied.  The log rank test compares the observed and expected number of events for each group using the same test statistic as the chi-squared test.



Survival Functions

---

Estimation of the Test Statistic for comparing two groups A and B: $\chi^2_{logrank} = \frac{(O_A - E_A)^2}{E_A} + \frac{(O_B - E_B)^2}{E_B}$

Where the expected number of events is calculated as: $E_{Aj} = \sum \frac{d_j n_{Aj}}{n_j}$

$d_j$ = no. of events at time $t_j$, $n_{Aj}$ = no. of people at risk at time j in group A and $n_j$ = total no. of people at risk.

Calculating the expected values is time consuming and the estimation of the test statistic is conservative. Instead you can use SPSS to calculate the test statistic and significance value.

Compare the test statistic with the critical value from the Chi-squared table.  The degrees of freedom are: number of groups – 1 so for a significance level of 5% and 1 df, the critical value is $\chi^2_1 = 3.84$.  If the test statistic is larger than 3.84, reject the null hypothesis and conclude that there is significant evidence to suggest a difference in survival times for the two groups.

Reviewer: Chris Knox
University of Sheffield

**Overall Comparisons**

|  | Chi-Square | df | Sig. |
|---|---|---|---|
| Log Rank (Mantel-Cox) | 3.971 | 1 | .046 |

Test of equality of survival distributions for the different levels of gender.

The p-value (sig) is the probability of getting a test statistic of at least 3.971 if there really is no difference in survival times for males and females.  As the p-value = 0.046 and is less than 0.05, conclude that there is significant evidence of a difference in survival times for males and females.  The estimated time until death is 2624 days for males and 1806 days for females this difference is statistically significant (p=0.046) therefore, males have an increased survival time compared to females following admission to hospital due to myocardial infarction.

**Cox's Regression**

Cox's regression allows several variables to be taken into account and tests the independent effects of these variables on the hazard of the event.

ANALYSE → SURVIVAL → COX REGRESSION



For all categorical variables, select the 'Categorical' option.  Tell SPSS whether you want to compare factor levels to the first or last category and select the 'Simple' option (although indicator appears to give the same main output).  Click on change after each variable.  Look at the help for info on other options but stick with simple for now.

For gender, selecting the reference category as first means that males are the reference category (as male = 0, female =1).

Hazard = the risk of reaching the event (e.g. death) at time point i, given that the individual has not reached it up to that time point, t. A lower hazard rate implies a higher survival rate. If the outcome is death the hazard rate can be interpreted as the mortality rate.

Model: $\lambda_i(t) = \lambda_0(t)\exp(\beta_1 x_1 + \beta_2 x_2 + ..... + \beta_k x_k)$  where:

There's a lot of output from SPSS but the following table contains the important output.

$\lambda_i(t)$ = the hazard function at time point t for individual I,

$\lambda_0(t)$ = the baseline hazard function (hazard function when all explanatory variables are set to 0)

**Variables in the Equation**

|  | B | SE | Wald | df | Sig. | Exp(B) | 95.0% CI for Exp(B) Lower | Upper |
|---|---|---|---|---|---|---|---|---|
| gender | .183 | .309 | .352 | 1 | .553 | 1.201 | .655 | 2.203 |
| bmi | -.071 | .036 | 3.860 | 1 | .049 | .931 | .867 | 1.000 |
| agegroup |  |  | 7.333 | 3 | .062 |  |  |  |
| agegroup(1) | .103 | .522 | .039 | 1 | .844 | 1.108 | .399 | 3.081 |
| agegroup(2) | .906 | .457 | 3.940 | 1 | .047 | 2.475 | 1.012 | 6.058 |
| agegroup(3) | 1.010 | .453 | 4.978 | 1 | .026 | 2.745 | 1.131 | 6.665 |

To understand the effects of individual predictors, look at Exp(β), which is the hazard ratio and can be interpreted as the predicted change in the hazard for a unit increase in the predictor.

---

**What's a Hazard Ratio?**

$$HR = \frac{h_A(t)}{h_B(t)} = \frac{risk\ of\ event\ (e.g.death)\ in\ group\ A}{risk\ of\ event\ in\ group\ B}$$

A hazard ratio can be interpreted in a similar way to relative risk.  It compares the risk of an event occurring in two groups.  If the ratio is above 1, the risk of the event happening in group A is higher.

---

 For gender, the reference category was males so the hazard ratio is $\frac{h_{females}(t)}{h_{males}(t)}$ = 1.201.  This indicates that the hazard (mortality) rate is 20% higher for females compared to males although the p-value and CI suggest this could be due to chance i.e. is non-significant.

The effect of BMI is statistically significant.  For each additional unit of BMI, the hazard decreases by $(1 - 0.931)*100$ = 6.9%.  For an additional 5 units increase in BMI, the hazard decreases by $(1 - 0.931^5)*100$ = 30%

For age group, the reference category is < 60.  The next category is 60 – 69, then 70 – 79 and lastly 80+.  For categorical variables, interpret Exp(β) directly e.g. the hazard (mortality) rate for 70 – 79 year olds is 147.5% higher than that of under 60's. The hazard (mortality) rate is 174.5% higher for the patients in the 80+ age group compared to those in the under 60's age group.

## Assumptions

There is one main assumption for survival analysis that is particularly important for Cox's regression. This is the proportional hazards assumption that the hazard ratio between two groups remains constant over time. Requesting a hazard plot in the cox regression menu gives you a figure like the one opposite. The lines should not cross each other and should be approximately parallel. If this is not the case modelling survival based on a distribution should be considered instead.

The proportional hazards assumption also applies to the log rank test and can be checked by assessing if the lines on the Kaplan-Meier plot remain parallel. If this assumption is violated the log-rank test has reduced power, in extreme cases it is an appropriate test to use.


Hazard Function for patterns 1 - 4

©Ellen Marshall
University of Sheffield

Reviewer: Chris Knox
University of Sheffield