

Cox proportional hazards regression

Professor Dr. Syed Hatim Noor
Unit of Biostatistics & Research Methodology
School of Medical Sciences
Universiti Sains Malaysia

Syed Hatim Noor

1

Cox proportional hazards regression

- To assess effect of multiple covariates on survival
- The most commonly used multivariable survival method
- Can be used to analyze data that contain censored observations

Syed Hatim Noor

2

Cox proportional hazards model

- Works with hazard model
- Can handle both continuous and categorical predictor variables
- **Proportional hazards assumption:**
 - Hazards ratio should be constant across time
- The proportionality of hazards should not vary over time

Syed Hatim Noor

3

Limitations of Cox PH model

- Covariates normally do not vary over time (**time independent covariates**)
 - True with respect to gender, ethnicity or congenital condition
- If proportional hazards assumption does not hold (hazard ratios change across time: value of one or more covariates are different at different time points)
 - Extended Cox regression model (**time-dependent or time-varying covariates**) needs to be applied

Syed Hatim Noor

4

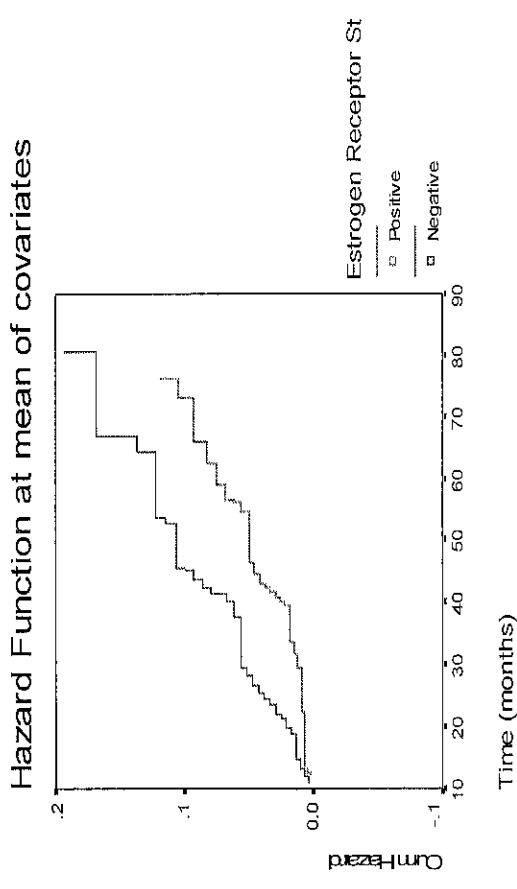
Hazard function $\{h(t)\}$

- Cox regression uses hazard function to estimate the hazards ratio
- A high hazard function indicates a high rate of mortality

Syed Hatim Noor

5

Hazard functions plot



Mathematical expression (Simple Cox)

- $h(t)=[h_0(t)] e^{bx}$
- $h(t)$ =hazard function (expected risk of death for a case with condition x)
- $h_0(t)$ =baseline hazard function (when x is set to zero) (expected risk without treatment or condition)
- b= regression coefficient
- x= independent prognostic variable
- e=base of the natural logarithm

Syed Hatim Noor

7

Explanation

- Dichotomous covariate x
- 0=control / no condition
- 1=treatment / condition
- $h(t)=[h_0(t)] e^{bx}$

Syed Hatim Noor

8

Expected risk of death with condition X

eg. Expected risk of death from breast cancer in a patient who has a positive estrogen receptor

$$h(t) = [h_0(t)] e^{bx}$$

↑

1=positive receptor status
0=negative receptor status

Expected risk of death with condition X

eg. Expected risk of death from breast cancer in a patient who has a negative estrogen receptor

Hazard ratio

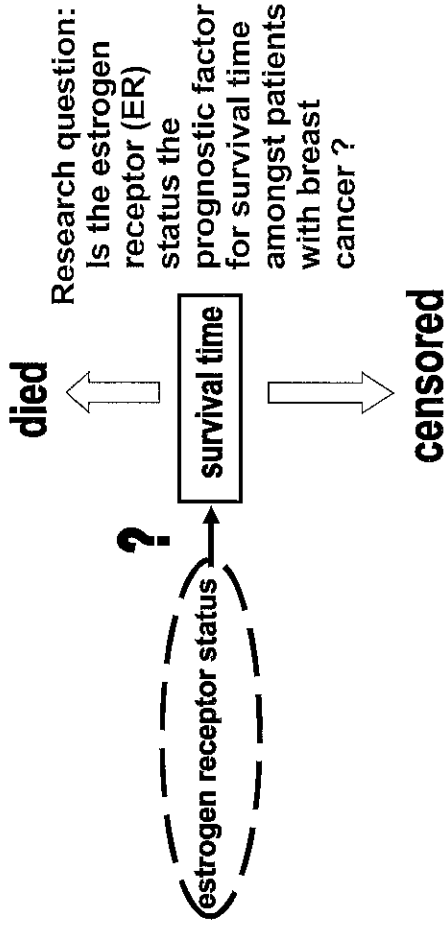
- Indicates the increase (or decrease) in risk incurred by applying the treatment or condition
- Also called as relative hazard
- $h(t)/[h_0(t)] = e^{bx}$
- Significance of the hazard ratios (HRs)
 - Confidence intervals
 - Statistical test

Steps in Cox proportional hazards regression model

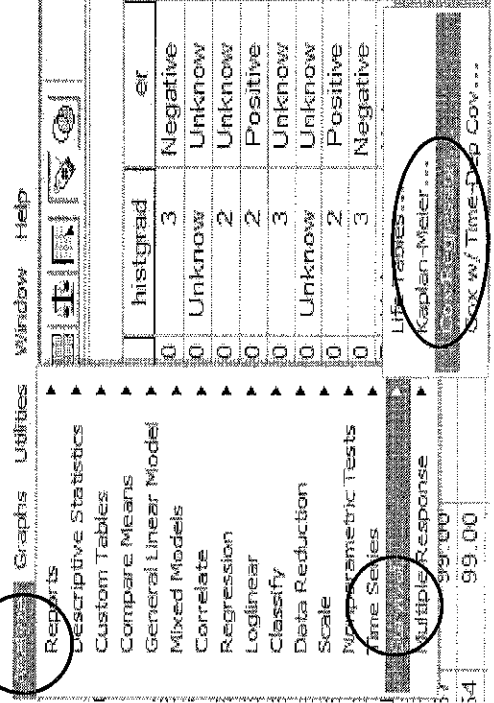
- Step 1: Data exploration and cleaning
- Step 2: Simple Cox regression / Kaplan-Meier survival analysis (log-rank test)
- Step 3: Variable selection (**prototyping**)
- Step 4: Checking multicollinearity and interactions (**fine modeling**)
- Step 5: Checking model assumptions
- Step 6: Interpretation, conclusion and presentation

Step 1: Data exploration and cleaning

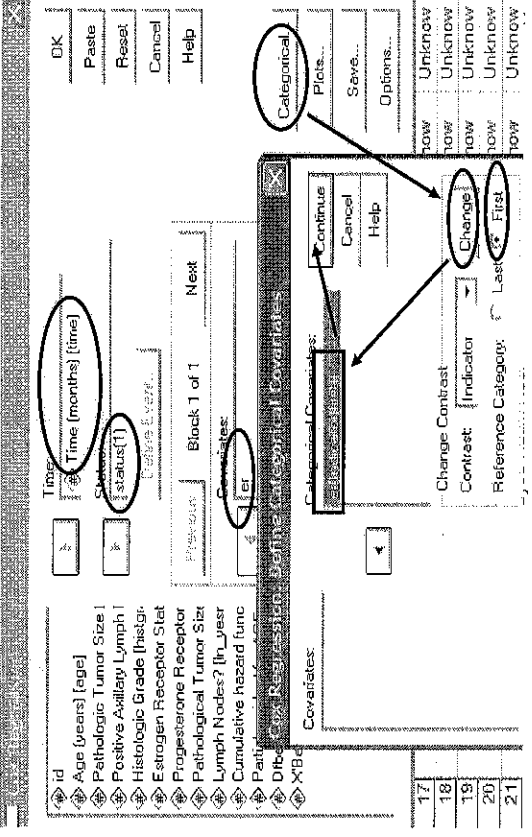
Example: breast cancer survival study (Simple Cox)



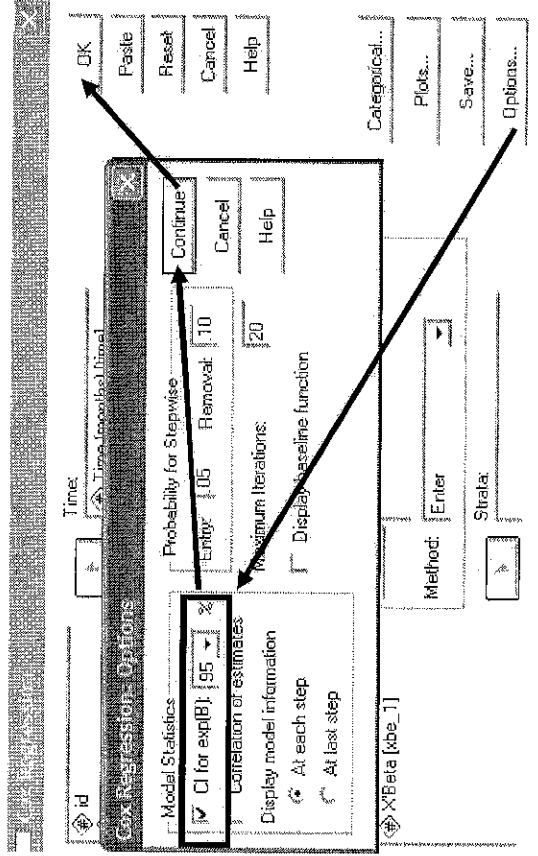
Steps in Simple Cox regression analysis



Steps in Simple Cox regression analysis



Steps in Simple Cox regression analysis



Hypothesis testing

- $H_0: B_1 = B_2 = B_3 = 0$
- Three main likelihood methods
 - Likelihood ratio test (LR test)
 - Wald test
 - Score test (global or overall chi-square test)

Case Processing Summary

	N	Percent
Cases available in analysis	51	4.2%
Censored	726	60.1%
Total	777	64.4%
Cases dropped		
Cases with missing value	338	28.0%
Cases with non-positive time	0	.0%
Censored cases before the event in a strat	92	7.6%
Total	430	35.6%
Total	1207	100.0%

a. Dependent Variable: Time (months)

Categorical Variable Codings

ER	Frequency	(1)
0=Negative	338	.000
1=Positive	531	1.000

a. Indicator Parameter Coding

b. Category variable: ER (Estrogen Receptor Site)

Omnibus Tests of Model Coefficients

	Overall (score)	Change From Previous Step	Change From Previous Block
-2 Log Likelihood	5.537		
Chi-square	5.537	5.5365	5.5365
df	1	1	1
Sig.	.019	.021	.021
Chi-square		5.365	5.365
df		1	1
Sig.		.021	.021

a. Beginning Block Number 0, initial Log Likelihood function: -2 Log likelihood: 613.088

b. Beginning Block Number 1, Method = Enter

$$613.088 - 607.723 = 5.365$$

- The negative B means that as the value of ER increases, the hazards decrease. coefficient is away from zero
- As the coding is 0=negative and 1=positive, hazards of dying is expected to be less when estrogen receptor status is positive.
- 95% confidence interval of hazards ratio does not include one and corresponding p-value is significant, estrogen receptor status appears to be the effective predictor variable when no other variables are considered.

25

Syed Hatim Noor

- The hazard function for estrogen receptor status can be written as

$$h(t) = [h_0(t)]e^{(-0.6507 * 1)}$$

- Since HR is less than one, it indicates there is a decreased hazards when estrogen receptor status is positive.

26

Syed Hatim Noor

Hands-on : Simple Cox

- Do Simple Cox regression analysis for all other independent variables (age, pathologic tumor size, number of positive lymph nodes, histologic grade, progesterone receptor status)
- Make a table for Simple Cox analysis results

27

Syed Hatim Noor

Step 3: Multiple Cox Proportional hazards Regression Variables selection (**prototyping**)

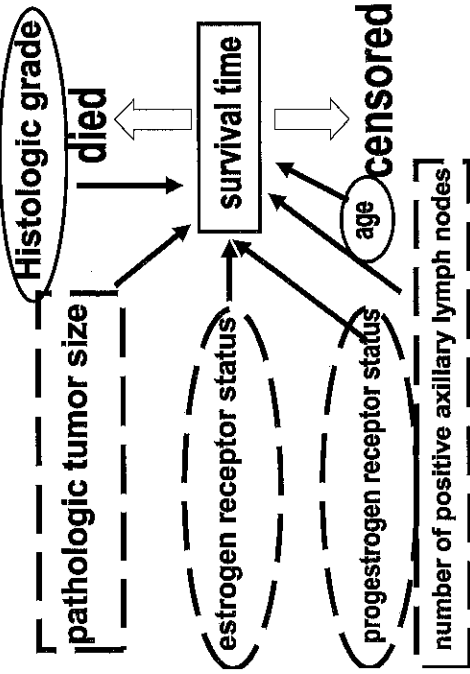
Syed Hatim Noor

28

Mathematical expression (Multiple Cox)

- $h(t) = [h_0(t)] e^{(B_1x_1 + B_2x_2 + B_3x_3 + \dots + B_nx_n)}$
- $h(t)$ = hazard function
- $h_0(t)$ = baseline hazard function
- B_1, B_2, B_3, \dots regression coefficients
- x_1, x_2, x_3, \dots prognostic factors (categorical / continuous)

Breast cancer survival study (Multiple Cox)

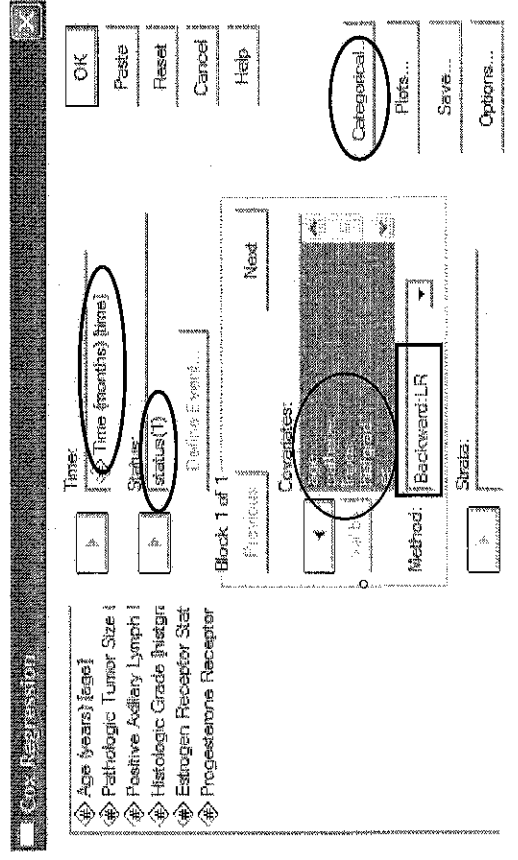


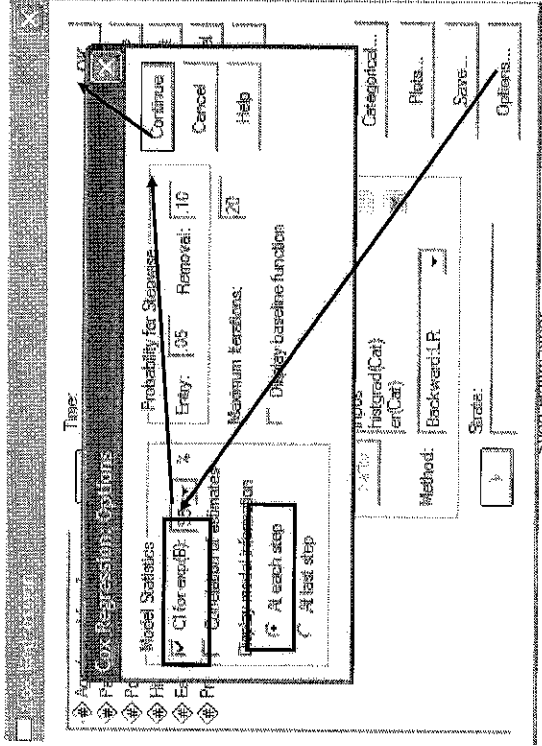
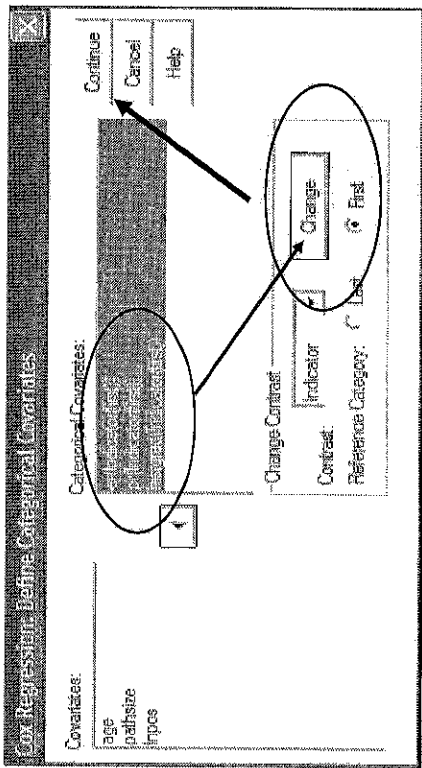
Research question:
Are estrogen receptor status, progesterone receptor status, number of positive axillary lymph nodes, pathologic tumor size, age and histologic grade the prognostic factors for survival time amongst patients with breast cancer?

Methods for selecting variables

- Forced entry (enter method) – all of the variables have been forced into the model in one step
- Forward, Backward and stepwise selections – variables are entered and deleted according to specified criteria

Steps in Multiple Cox Regression





Case Processing Summary

	N	Percent
Cases available in analysis	40	3.8%
Censored	550	45.6%
Total	590	48.9%
Cases dropped	347	45.3%
Cases with missing values	0	.0%
Cases with negative Censored cases before the earliest event in stratum	70	5.8%
Total	617	51.1%
Total	1207	100.0%

a. Dependent Variable: Time (months)

Categorical Variable Coding^a

	Frequency	(1)	(2)
histgrad 1=1	56	0	0
2=2	352	1	0
3=3	252	0	1
e ^a 0=Negative	262	0	0
1=Positive	398	1	1
p ^a 0=Negative	299	0	0
1=Positive	361	1	1

a. Indicator Parameter Coding

b. Category variable: histgrad (Histologic Grade)

c. Category variable: er (Estrogen Receptor Status)

d. Category variable: pr (Progesterone Receptor S

Variables in the Equation

Step	Variable	B	SE	Wald	df	Sig.	Exp(B)	95.0% CI for Exp(B)	
								Lower	Upper
Step 1	age	-.020	.014	2.007	1	.157	.980	.954	1.008
	pathsize	.392	.140	7.894	1	.005	1.480	1.126	1.946
	Inpos	.133	.049	7.342	1	.007	1.143	1.038	1.258
Step 2	histgrad	.866	1.033	.427	2	.658			
	histgrad(1)	.675	1.033	.427	1	.513	1.965	.259	14.890
	histgrad(2)	.866	1.049	.682	1	.409	2.378	.304	18.596
Step 3	er	.078	.431	.083	1	.857	1.081	.464	2.518
	pr	-.542	.425	1.629	1	.202	.581	.253	1.337
	age	-.020	.014	1.980	1	.159	.981	.954	1.008
Step 4	pathsize	.389	.138	7.904	1	.005	1.475	1.125	1.934
	Inpos	.134	.049	7.382	1	.007	1.143	1.038	1.259
	histgrad	.805	.805	.805	2	.669			
Step 3	histgrad(1)	.674	1.033	.426	1	.514	1.963	.259	14.876
	histgrad(2)	.852	1.046	.664	1	.415	2.345	.302	18.239
	pr	-.499	.353	2.002	1	.157	.607	.304	1.212
Step 3	age	-.022	.014	2.525	1	.112	.979	.953	1.005
	pathsize	.395	.138	8.217	1	.004	1.484	1.133	1.944
	Inpos	.138	.049	7.965	1	.005	1.148	1.043	1.254
Step 4	pr	-.561	.344	2.650	1	.104	.571	.291	1.121
	pathsize	.449	.132	11.565	1	.001	1.566	1.209	2.028
	Inpos	.137	.047	8.369	1	.004	1.147	1.045	1.259
Step 4	pr	-.672	.336	4.001	1	.045	.511	.265	.987

Results and Interpretation

- Each unit increase in pathologic tumor size is expected to increase the hazards of dying by 1.6 times (b=0.45, HR=1.57, 95% CI 1.21,2.03, p=0.001)
- Each unit increase in number of positive lymph nodes is expected to increase the hazards of dying by 1.2 times (b=0.14, HR=1.15, 95% CI 1.05,1.26, p=0.004)
- Hazards of dying is decreased by half when progesterone receptor is positive (b=-0.67, HR=0.51, 95% CI 0.27, 0.99, p=0.045)

$$h(t)=[h_0(t)] e^{(0.45 \times \text{pathsize}) + (0.14 \times \text{lnpos}) + (-0.67 \times \text{pr})}$$

Step	pathsize	.449	.122	11.565	1	.001	1.566	1.209	2.028
4	lnpos	.137	.047	8.368	1	.004	1.147	1.045	1.259
	pr	-.672	.336	4.001	1	.045	.511	.285	.987

Syed Hatim Noor

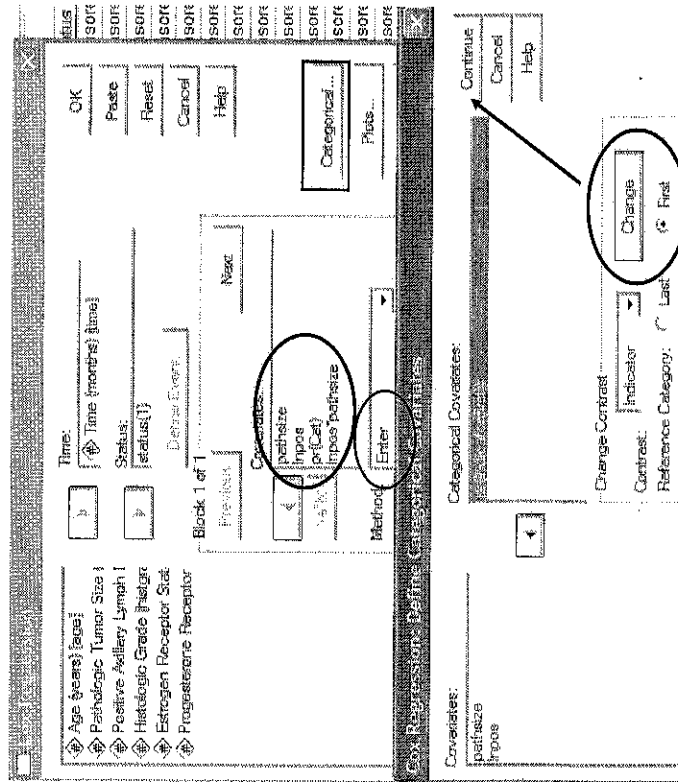
Model if Term Removed

Term Removed	Loss	Chi-square	df	Sig.
Step 1	age	2.054	1	.152
	pathsize	7.183	1	.007
	lnpos	6.044	1	.014
	histgrad	.996	2	.617
	er	.033	1	.857
	pr	1.636	1	.201
Step 2	age	2.022	1	.155
	pathsize	7.165	1	.007
	lnpos	6.080	1	.014
	histgrad	.934	2	.627
	pr	2.050	1	.152
Step 3	age	2.598	1	.108
	pathsize	7.446	1	.006
	lnpos	6.528	1	.011
	pr	2.727	1	.099
Step 4	pathsize	10.251	1	.001
	lnpos	6.778	1	.009
	pr	4.158	1	.041

a. Residual Chi Square = .033 with 1 df Sig. = .857
 b. Residual Chi Square = .864 with 3 df Sig. = .834
 c. Residual Chi Square = 3.334 with 4 df Sig. = .50

Syed Hatim Noor

Step 4: Checking interactions (fine modeling)



Syed Hatim Noor

Variables in the Equation

	B	SE	Wald	df	Sig.	Exp(B)	95.0% CI for Exp(B)	
							Lower	Upper
pa1nsize	-.562	.141	15.891	1	.000	1.754	1.331	2.312
lnpos	-.263	.074	12.798	1	.000	1.301	1.127	1.504
pr	-.742	.307	5.865	1	.015	.476	.261	.868
lnpos*pa1nsize	-.044	.027	2.692	1	.101	.957	.908	1.009

Variables in the Equation

	B	SE	Wald	df	Sig.	Exp(B)	95.0% CI for Exp(B)	
							Lower	Upper
pa1nsize	-.367	.132	8.585	1	.003	1.473	1.137	1.908
lnpos	.149	.037	16.541	1	.000	1.161	1.080	1.247
pr	-.818	.660	1.535	1	.215	.441	.121	1.610
pa1nsize*pr	.070	.252	.078	1	.780	1.073	.656	1.758

Variables in the Equation

	B	SE	Wald	df	Sig.	Exp(B)	95.0% CI for Exp(B)	
							Lower	Upper
pa1nsize	.446	.125	12.694	1	.000	1.562	1.222	1.998
lnpos	.105	.085	2.574	1	.109	1.111	.977	1.262
pr	-.788	.343	5.261	1	.022	.455	.232	.892
lnpos*pr	-.062	.076	.667	1	.414	1.064	.917	1.234

Syed Hatim Noor

Step 5: Checking model assumptions

Syed Hatim Noor

Graphical approaches

- Several graphical approaches
 - Check proportionality assumption with **hazard functions plot and Log-minus-log plots** of the baseline survival functions
 - Check residuals
 - Partial residual (Schoenfeld residuals)**
 - Martingale residuals**
 - Cox-Snell residuals (cumulative hazard function)**
 - Deviance residuals**
 - X'Beta**
 - Influential statistics**
 - DfBeta**

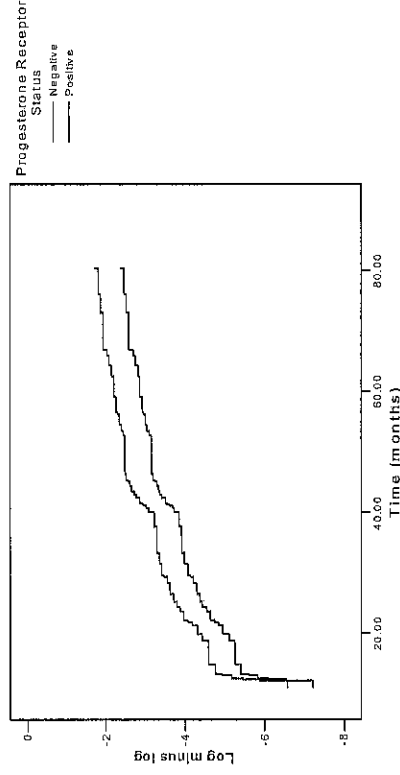
Syed Hatim Noor

- Proportional assumption: covariates are independent with respect to time and their hazards ratio are constant over time
- General ways to examine model adequacy
 - Graphically
 - Mathematically

Syed Hatim Noor

Log-minus-log (LML) plot

LML Function for patterns 1 - 2



Interpretation : The curves seem sufficiently parallel. The proportional hazard assumption is met.

Syed Hatim Noor

49

Syed Hatim Noor

50

Checking proportional hazards assumption

- Can be checked by examining the (LML) log-minus-log survival plot
- Should display parallel lines
- Can be checked only for categorical variables

Checking residuals

- ⇒ Partial residual (Schoenfeld residuals)
- Martingale residuals
- Cox-Snell residuals (cumulative hazard function)
- Deviance residuals
- X'Beta

Syed Hatim Noor

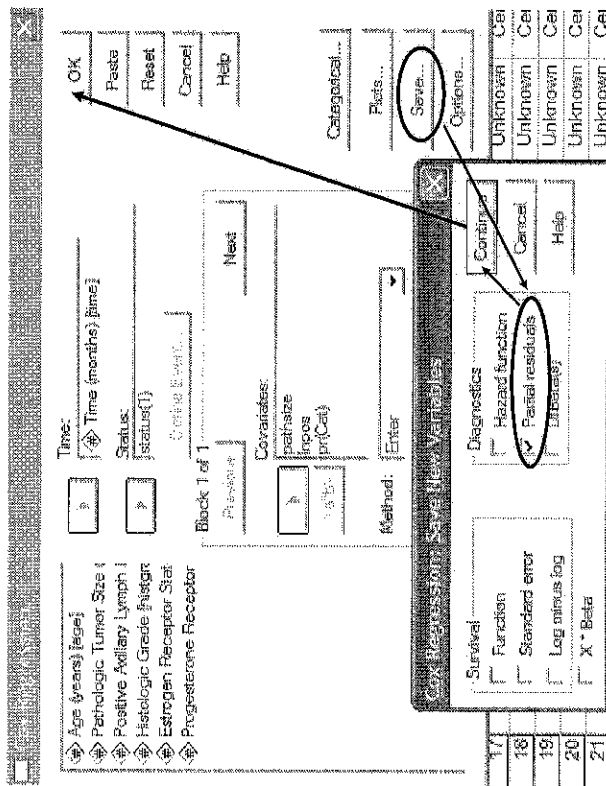
51

Partial residuals (Schoenfeld residuals)

- Can be plotted against time to test for violations of the proportional hazards assumption
- Partial residual for each case is the difference between the observed value of the case and its expected value
- Difference should be approximately zero if proportionality assumption holds

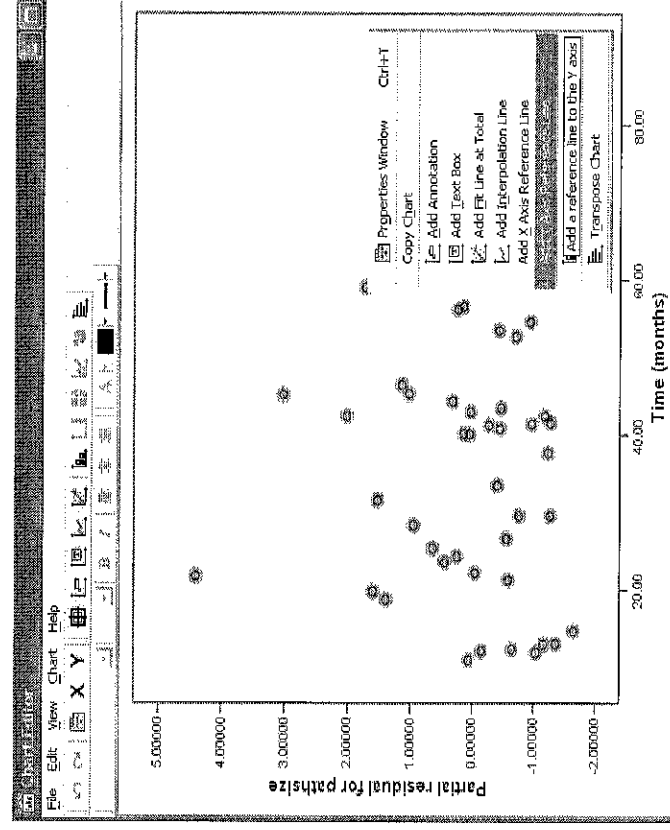
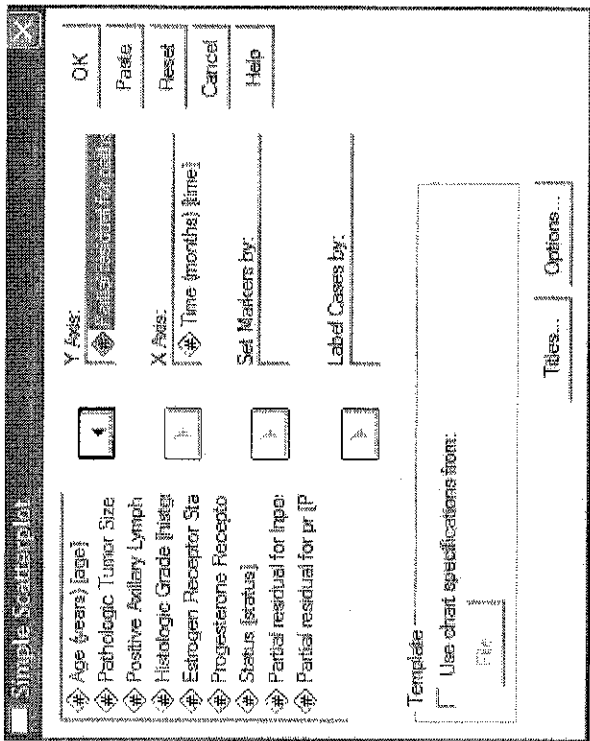
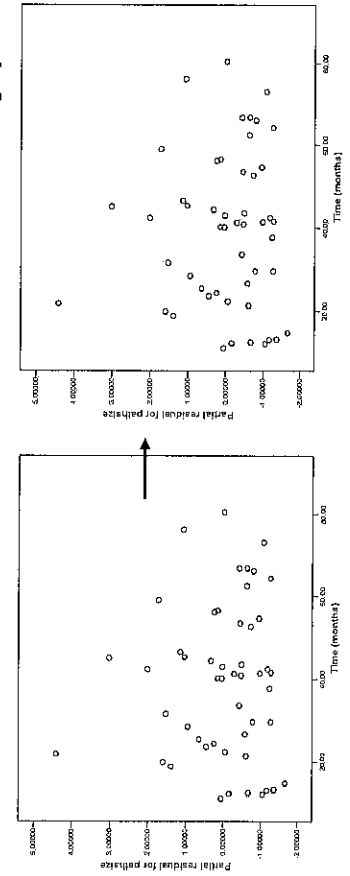
Syed Hatim Noor

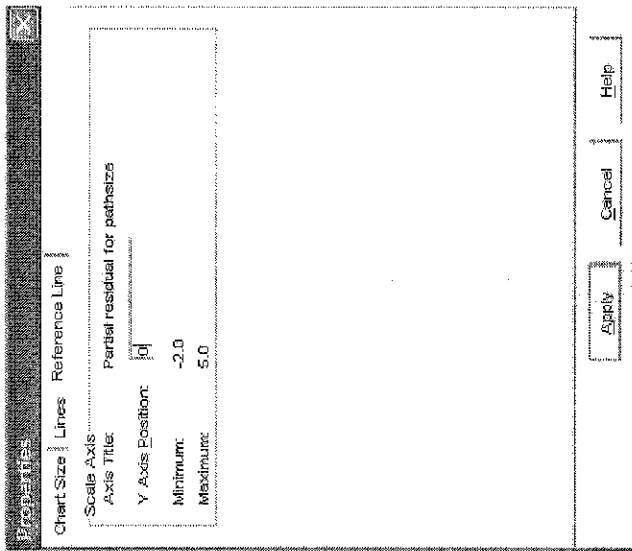
52



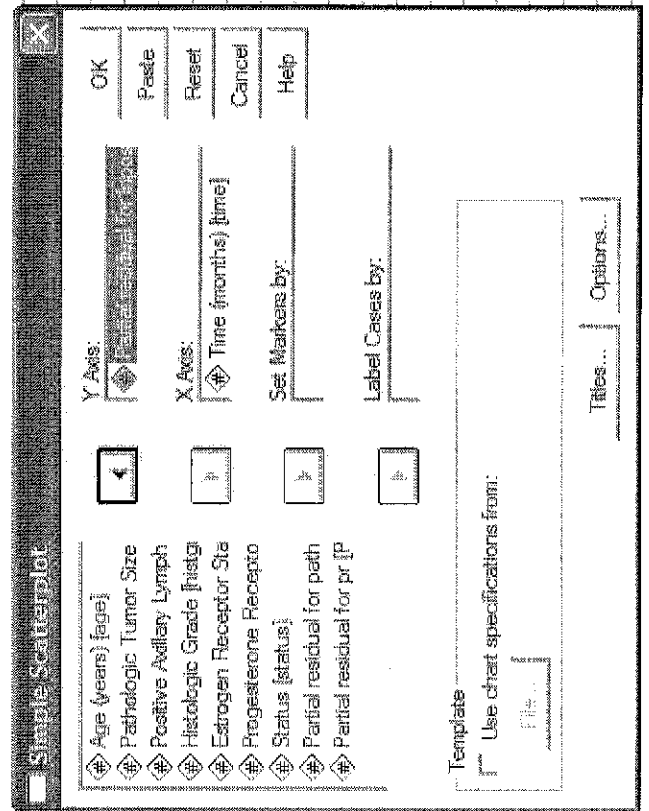
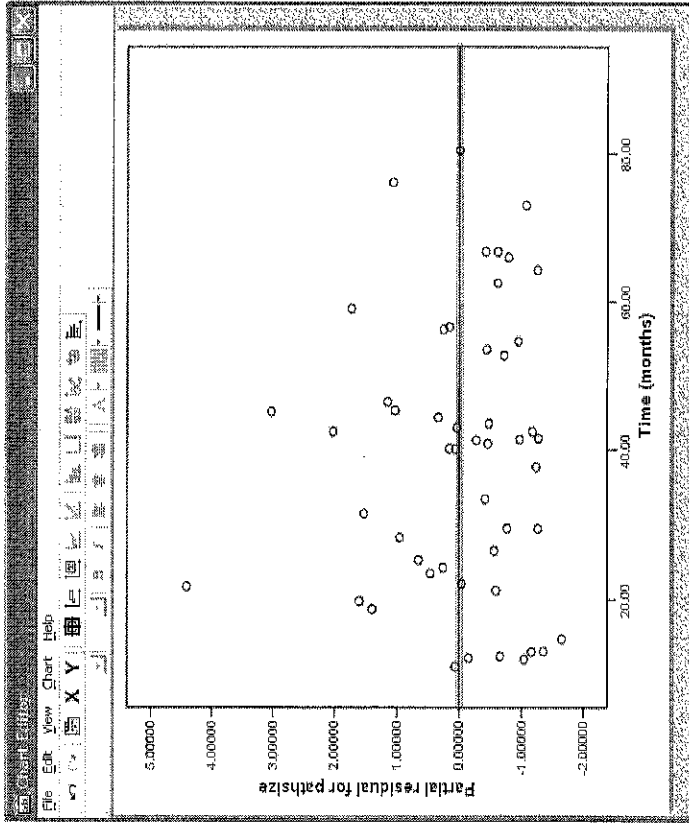
Steps in getting a partial residuals plot

double click on the graph





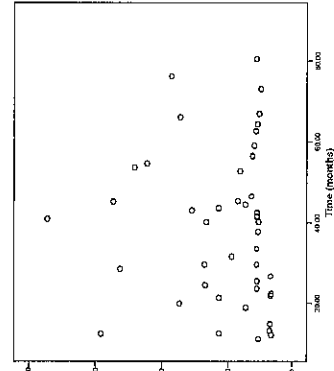
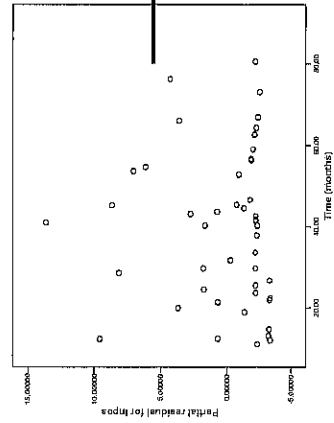
57

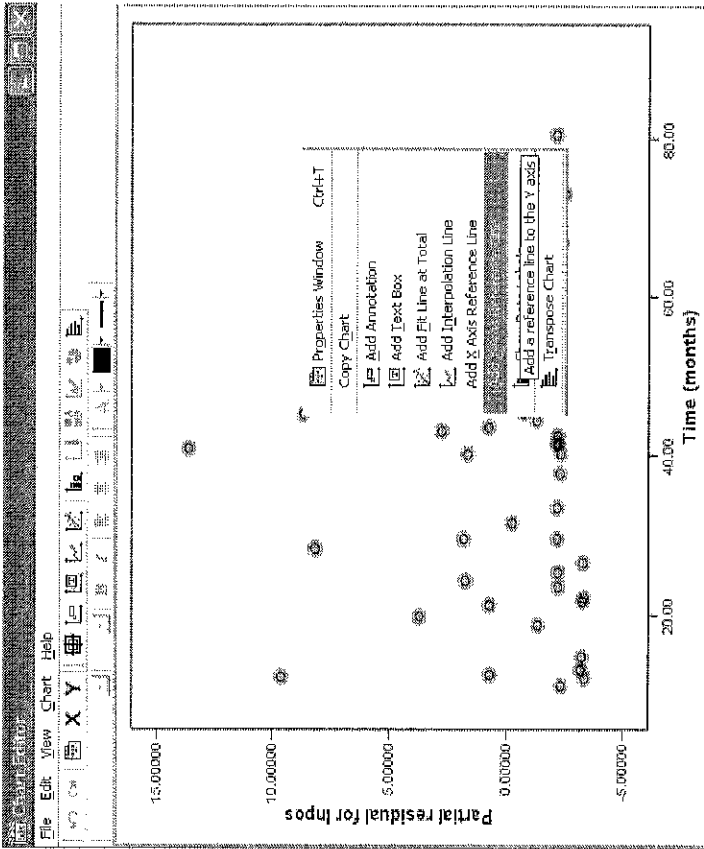


59

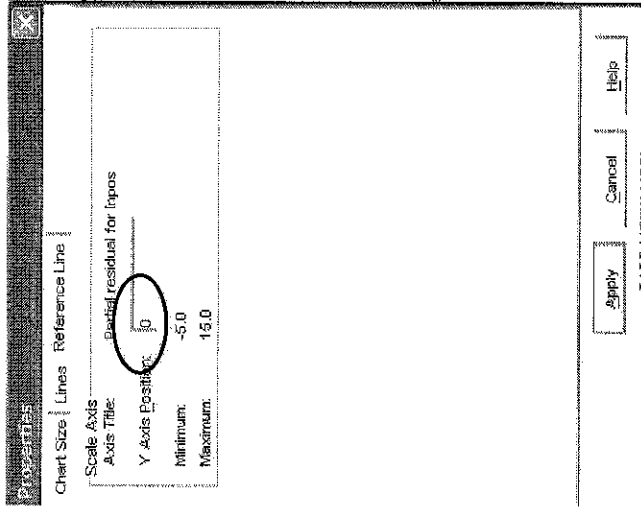
Steps in getting a partial residuals plot

double click on the graph

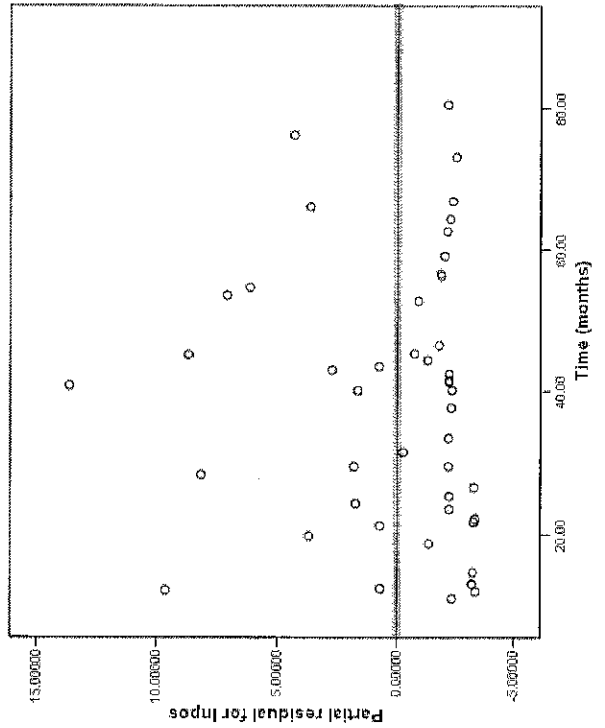




61



62

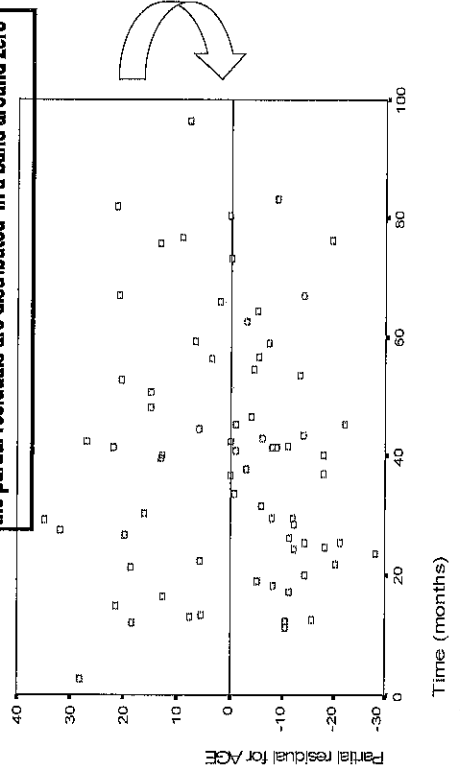


Syed Hatim Noor

63

Other example

the partial residuals are distributed in a band around zero



Assumption of proportional hazards is met

64

Hands-on Exercise

By using the data set (breast cancer), perform the following data analyses:

- (1) Use the Simple Cox regression model to determine whether the variables (estrogen receptor status, progesterone receptor status, age, pathological tumor size, number of positive lymph nodes and histologic grade) are potential prognostic factors of survival time of patients with carcinoma of the breast. Comment on the findings.

Step 6:

Interpretation, presentation and conclusion

Hands-on exercise

- (2) Perform multiple Cox regression to determine the prognostic factors of breast cancer by applying different methods for selecting models (forced entry, forward stepwise and backward stepwise). Compare the results and make comments.
- (3) Check possible clinically sound interactions between covariates by including two-way interactions terms.

Hands-on exercise

- (4) Check whether the model holds the proportional hazards assumption by applying hazard function plot and Log-minus-log plot. Check each covariate from the final model with regard to partial residuals (Schoenfeld residuals) to check whether proportional hazards assumption is met.

Hands-on exercise

- (5) Interpret the results and make the conclusion.
- (6) Present the results in tables.

Summary

- Survival analyses quantifies **time to a dichotomous event**
- Handles **censored data** well
- **Kaplan-Meier survival analysis** can be compared statistically and graphically
- **Cox proportional hazards models** help distinguish individual contributions of covariates on survival, provided certain assumptions are met

Table (1) Prognostic factors of cancer breast by Simple Cox proportional hazards model

Variable	Regression coefficient (b)	Crude Hazards ratio (95% CI)	Wald statistic	p-value
Age	-0.02	0.98(0.96,0.99)	4.98	0.026
Pathologic tumor size	0.60	1.82(1.53,2.17)	45.28	<0.001
Number of positive lymph nodes	0.10	1.11(1.07,1.15)	29.47	<0.001
Histologic grade				
Grade 1	0.00	1.00	-	-
Grade 2	0.71	2.03(0.48,8.55)	0.94	0.333
Grade 3	1.35	3.86(0.92,16.25)	3.40	0.065
Estrogen receptor status				
Negative	0.00	1.00	-	-
positive	-0.65	0.52(0.30,0.91)	5.35	0.021
Progesterone receptor status				
Negative	0.00	1.00	-	-
positive	-0.64	0.53(0.30,0.93)	4.97	0.026

Table (2) Prognostic factors of cancer breast by multiple Cox proportional hazards model

Variable	Regression coefficient (b)	Adjusted Hazards ratio (95% CI)	Wald statistic	p-value
Pathologic tumor size	0.45	1.57(1.21, 2.03)	11.57	0.001
Number of positive lymph nodes	0.14	1.15(1.05, 1.26)	8.37	0.004
Progesterone receptor status				
Negative	0.00	1.00	-	-
positive	-0.67	0.51(0.27, 0.99)	4.00	0.045

Backward stepwise Cox proportional hazards regression model applied
 Log-minus-log plot, hazard function plot and partial residuals were applied to check the model assumption and found fulfilled

Table (3) Prognostic factors of cancer breast by univariable and multivariable Cox proportional hazards models

Variable	Simple Cox regression		Multiple Cox regression			
	b	Crude HR(95%CI)	p	b	Adjusted HR(95% CI)	p
Number of positive lymph nodes	0.10	1.11(1.07,1.15)	<0.001	0.14	1.15(1.05,1.26)	0.004
Progesterone receptor status						
Negative	0.00	1.00	-	0.00	1.00	-
positive	-0.64	0.53(0.30,0.93)	0.026	-0.67	0.51(0.27, 0.99)	0.045
Pathologic tumor size	0.60	1.82(1.53,2.17)	<0.001	0.45	1.57(1.21, 2.03)	0.001

Backward stepwise Cox proportional hazards regression model applied
 Log-minus-log plot, hazard function plot and partial residuals were applied to check the model assumption and found fulfilled