



Australian
National
University

RESEARCH SCHOOL OF FINANCE, ACTUARIAL STUDIES
AND APPLIED STATISTICS

First Semester Final Examination (2014)

Survival Models / Biostatistics
(STAT3032/7042/8003)

Writing period: 3 hours duration

Study period: 15 minutes duration

*Permitted materials: Non-programmable calculator, dictionary,
one A4 sized sheet of paper with notes on both sides*

Total marks: 70 (undergraduates) / 78 (postgraduates)

INSTRUCTIONS TO CANDIDATES:

- *Postgraduates should attempt all questions. Undergraduates should only attempt questions 1 to 7.*
- *To ensure full marks show all the steps in working out your solutions. Marks may be deducted for failure to show appropriate calculations or formulae.*
- *All questions are to be completed in the script book provided.*
- *All answers should be rounded to 4 decimal places.*

Question 1 [7 marks]

For a particular population it is shown that $l_x = 100 - x$, $0 \leq x \leq 100$. Using this information about the number of lives aged x exact, calculate the following:

- (a) [2 marks] the force of mortality at age 50.
- (b) [2 marks] the complete expectation of life at age 50.
- (c) [3 marks] the average age of individuals who die between ages 40 and 45.

Question 2 [8 marks]

(For each part, you will gain 2 marks for a correct answer, be penalized 1 mark for an incorrect answer, and score 0 if no answer is given.) Answer each question “TRUE” or “FALSE”. In each case, write the whole word. It is **not** acceptable to write only “T” or “F” and answers presented in this form **will be graded incorrect**.

- (a) [2 marks] The force of mortality function μ_x can be greater than 1. ✓
- (b) [2 marks] If the force of mortality function μ_x is assumed to follow Makehams law, this means that for each one year increment in age, μ_x increases by a constant additive amount. ✗
- (c) [2 marks] Increasing the bandwidth of a Kernel function will lead to a rougher smoothed curve. ✗
- (d) [2 marks] Simultaneously conducting n ($n > 1$) multiple independent hypothesis tests (all at α % significance level) will lead to a new significance level which is theoretically equal to α^n %. ✗

Question 3 [9 marks]

In this question you will be looking at the breastfeeding data that was discussed in class. As a reminder, this dataset contains information on the duration of breastfeeding for over 900 mothers as well as information on a number of covariates. To investigate the effects of these covariates a Cox regression model was fitted. The following is output from a Cox regression analysis conducted in R:

	coef	exp(coef)	se(coef)	z	Pr(> z)	
as.factor(race)2	0.187823	1.206620	0.105469	1.781	0.074940	.
as.factor(race)3	0.296366	1.344962	0.097186	3.049	0.002293	**
as.factor(poverty)1	-0.226792	0.797087	0.094469	-2.401	0.016363	*
as.factor(smoke)1	0.246100	1.279027	0.079599	3.092	0.001990	**
as.factor(alc)1	0.167771	1.182666	0.123253	1.361	0.173454	
agemth	-0.173193	0.840975	0.186168	-0.930	0.352212	
I(agemth^2)	0.003615	1.003621	0.004251	0.850	0.395118	
ybirth	0.079672	1.082932	0.020505	3.886	0.000102	***
yschool	-0.056739	0.944841	0.023167	-2.449	0.014320	*

The variables *race*, *poverty*, *smoke*, and *alcohol* are categorical variables representing the race of the mother, whether the mother is living in poverty, whether the mother smokes, and whether the mother consumes alcohol. The variables *agemth*, *ybirth*, and *yschool* are continuous variables representing the age of the mother, the year of birth of the child and the number of years that the mother attended school. Based on the above output answer the following questions:

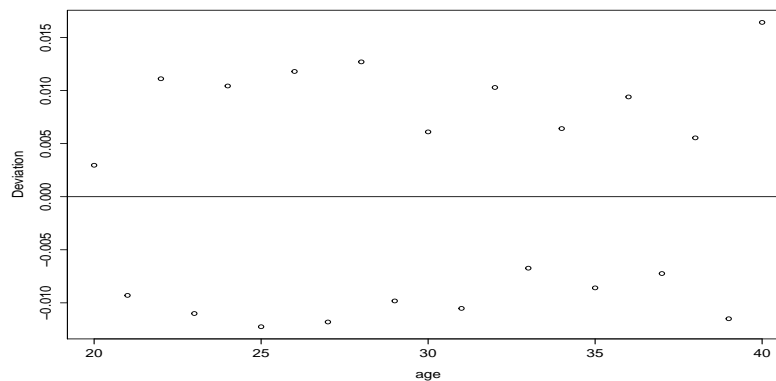
- [2 marks]** Use a critical value of 2 to provide a 95% confidence interval for the multiplicative change in the hazard for a 1 year increase in the variable *yschool*, everything else held constant.
- [3 marks]** Provide an estimate of the quantity $\frac{1}{\beta_9}$, where β_9 is the parameter corresponding to *yschool*. You must provide a standard error for your estimate.
- [2 marks]** After the above analysis had been conducted you are told that the values of duration were mistakenly reported using units of months rather than days. For example, a duration of 4 days was accidentally recorded as a duration of 4 months. How would this information change the parameter estimates reported above? You must provide a reason for your answer.
- [2 marks]** What would happen to the value of the partial likelihood if the coefficient estimate for the variable *yschool* was replaced by the value 0.05 (note all other estimates remain unchanged). You must provide a reason for your answer.

Question 4 [16 marks]

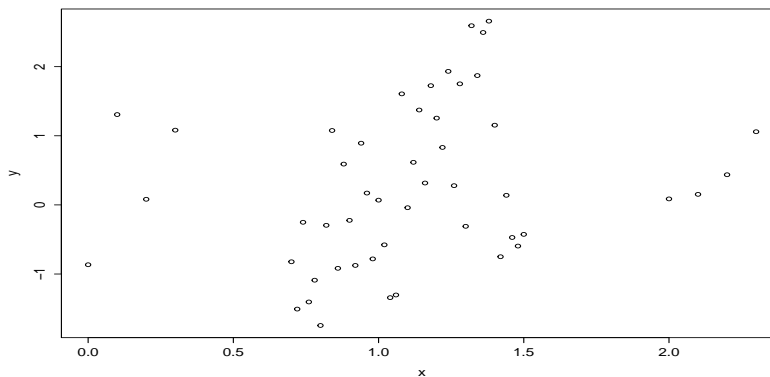
A set of crude mortality rates were graduated using a particular smoothing technique. **The smoothing technique that was used required 5 parameters to be estimated.** The table below provides details of the graduation:

Age	Crude Rate	Graduated Rate	Exposed to Risk	Deaths	Expected Deaths	Standardised Deviation	Squared Std. Dev.
20	0.028884	0.02287	1004	29	22.96	1.2748	1.6251
21	0.016260	0.02525	984	16	24.85	-1.7975	3.2310
22	0.040123	0.02778	972	39	27.00	2.3416	5.4831
23	0.018386	0.03049	979	18	29.85	-2.2027	4.8519
24	0.043750	0.03339	960	42	32.05	1.7867	3.1923
25	0.022293	0.03648	942	21	34.36	-2.3225	5.3940
26	0.052295	0.03978	937	49	37.27	1.9601	3.8420
27	0.031083	0.04332	933	29	40.42	-1.8361	3.3713
28	0.061159	0.04712	932	57	43.92	2.0226	4.0909
29	0.039474	0.05122	912	36	46.71	-1.6092	2.5895
30	0.060241	0.05566	913	55	50.82	0.6037	0.3645
31	0.048206	0.06047	892	43	53.94	-1.5367	2.3614
32	0.078409	0.06570	880	69	57.82	1.5217	2.3156
33	0.067045	0.07139	880	59	62.82	-0.5006	0.2506
34	0.087558	0.07759	868	76	67.35	1.0977	1.2049
35	0.077816	0.08434	861	67	72.62	-0.6888	0.4744
36	0.104094	0.09167	855	89	78.38	1.2589	1.5848
37	0.092216	0.09963	835	77	83.19	-0.7153	0.5117
38	0.112961	0.10824	841	95	91.03	0.4406	0.1941
39	0.102871	0.11752	836	86	98.25	-1.3153	1.7300
40	0.142327	0.12747	808	115	103.00	1.2663	1.6035
Total	1.327451	1.31738	19024	1167	1158.61	1.0500	50.2670

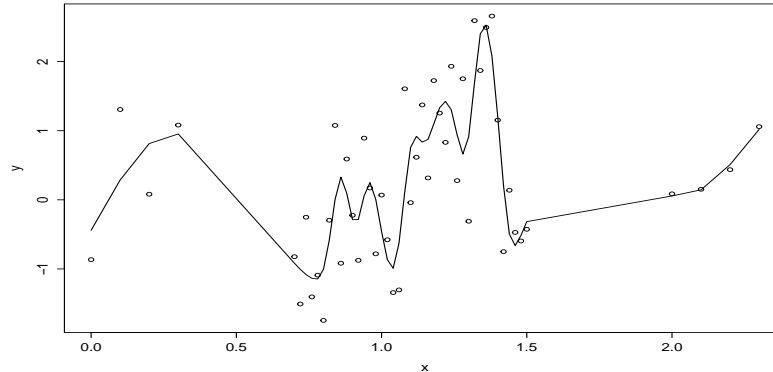
- (a) **[3 marks]** Perform the chi-square test to check whether the graduated rates are appropriate. State the required steps discussed in the tutorial (the test statistic, critical value and conclusion of the statistical test). Use a significance level of 5%.
- (b) **[3 marks]** Perform the standardized deviation test with four regions $(-\infty, -1.5]$, $(-1.5, 0]$, $(0, 1.5]$ and $(1.5, \infty)$. State the required steps as in part (a) with a significance level of 5%.
- (c) **[3 marks]** Consider the ages 20-25 only, calculate the kernel smoothed rate for age 22 from the crude rates provided in the table, using the triangle kernel with bandwidth 5.
- (d) **[3 marks]** The figure below shows the deviations resulting from the triangle kernel based on the entire sample. Using this figure to perform the sign test. State the required steps as in part (a) with a significance level of 5%.



- (e) **[2 marks]** The figure below shows a plot of a particular set of data that is going to be smoothed using kernel smoothing (in conjunction with the Normal kernel). Would you have any concerns smoothing the data depicted in this figure using kernel smoothing?



- (f) [2 marks] We use the natural cubic spline to generate the smoothed curve shown below. Do you think this is an appropriate choice of knots? If not, how would you like to adjust the number of knots?



Question 5 [10 marks]

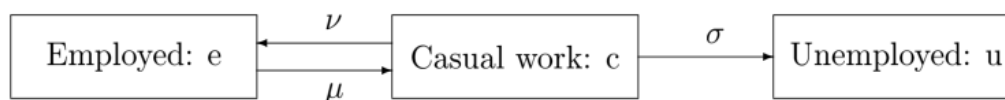
The table gives data on 18 lives from a portfolio of policies on impaired lives. The lives are classified into two categories, Smoking and Non-smoking. The table gives the time in months until a claim is made; the * indicates that the observation was censored.

Smoking	2	7*	9	10*	10*	10*	11	14	17
Non-smoking	4*	5	10*	10*	10*	12*	13	15*	18*

- (a) [4 marks] Find the Nelson-Aalen estimate of the survivor function for all lives combined. *Note: You don't need to calculate the standard errors.*
- (b) [2 marks] In a Cox regression $\mu_x(t) = \mu_0(t)e^{\beta x}$, where $x = 0$ for smoking and $x = 1$ for non-smoking, show the **partial log likelihood function** $l(\beta)$. Simplify your results as much as possible. *Note: **DO NOT** ignore any terms of $l(\beta)$.*
- (c) [4 marks] Apply the Maximum Likelihood method to estimate β . In the calculation, use the approximation $a + (a+1)e^\beta \approx (a+1)(1+e^\beta)$. Provide the corresponding standard error. Use Z-test to test whether $\beta = 0$, you only need to **report the p-value**. Base your conclusion on 5 % significance level.

Question 6 [12 marks]

The following three state Markov model with constant transition intensities μ , ν and σ is used to represent the transition from employment to long-term unemployment.



- (a) [6 marks] Show that ${}_t p_x^{eu}$ satisfies the differential equation

$$\frac{\partial^2}{\partial t^2} {}_t p_x^{eu} + (\mu + \sigma + \nu) \frac{\partial}{\partial t} {}_t p_x^{eu} + \sigma \mu {}_t p_x^{eu} = \sigma \mu$$

[Hint: show the first order of derivatives of ${}_t p_x^{eu}$ and ${}_t p_x^{ec}$ and then consider ${}_t p_x^{eu} + {}_t p_x^{ee} + {}_t p_x^{ec}$.]

Now only consider the two-state case (the Casual work and Unemployed). Suppose we have the below observations ($\delta = 1$ for death and $\delta = 0$ for censoring) for a group of people aged x and the force of mortality is a constant (σ) in $[x, x + 1]$:

x	a_i	b_i	δ_i	T_i
1	0	1	0	1
2	0.4	0.8	0	0.8
3	0.5	0.7	1	0.6
4	0	1	1	0.8

- (b) [2 marks] Using the binomial model of mortality obtain the exact likelihood arising from the above data. Simplify your results as much as possible.
- (c) [4 marks] Using the two-state model of mortality to calculate the maximum likelihood estimate of σ , and hence the maximum likelihood estimate of q_x . Provide a standard error of \hat{q}_x .

Question 7 [8 marks]

A life office has carried out an investigation of the mortality experience of certain of their policyholders. Part of the data is given below.

Age	Census Data					Total Deaths
	31/12/95	30/6/96	31/12/96	30/6/97	31/12/97	
50	5451	4420	4515	4773	5934	70
51	6002	5722	5534	4631	4428	76
52	5789	6121	6087	5322	5172	83

- (a) **[3 marks]** Suppose first that the definitions of age for the census data and death data are the same. Estimate the central exposed to risk E_{51}^c . Hence estimate μ_{51} and q_{51} . When calculating q_{51} , assuming that μ_{51} is constant and **DO NOT** use the initial exposure to risk.
- (b) **[3 marks]** Suppose now that the definition of age for the census data is: age x nearest birthday; the definition of age for those who died is: age x next birthday. Use this additional information to re-estimate E_{51}^c , μ_{51} and q_{51} . When calculating q_{51} , assuming that μ_{51} is constant and **DO NOT** use the initial exposure to risk.
- (c) **[2 marks]** Using the poisson model to calculate an approximate standard error for $\hat{\mu}_{51}$ derived in part (b).

Question 8 **[8 marks]** (For students enrolled in STAT7042/8003 ONLY)

(For each part, you will gain 2 marks for a correct answer, be penalized 1 mark for an incorrect answer, and score 0 if no answer is given.) Answer each question “TRUE” or “FALSE”. In each case, write the whole word. It is **not** acceptable to write only “T” or “F” and answers presented in this form **will be graded incorrect**.

- (a) **[2 marks]** Lasso regression can lead to exact 0 for some estimated parameters.
- (b) **[2 marks]** Estimates obtained from Ridge regression are unbiased.
- (c) **[2 marks]** We cannot use the Bootstrap to calculate confidence interval for the median.
- (d) **[2 marks]** The test statistic r_j of the j th Serial Correlation asymptotically follows the standard normal distribution.

END OF EXAMINATION