## Statistics 520 Midterm Exam I

Sept. 22, 2023

Yi = survival states of $x = 1000$   Same space of $x = 1000$   Some space			50	P** ==,				
INSTRUCTIONS: Read the questions carefully and completely. Answer each question and show all your work in the space provided. Credit cannot be given if work is not shown. Good luck!  Analysis of Survival Data for Cancer Treatment (25 pointes).  A local hospital is analyzing the effectiveness of two different cancer treatments: A and B. The hospital wants to assess whether the survival rate six months after treatment is significantly different between the two treatments. Treatment is considered effective if a patient survives for at least six months after receiving it.  Data has been collected from 12 patients who received treatment, with survival after 6 months recorded as 1 (Survived) or 0 (Did not survive). 6 patients are randomly selected to receive treatment A, and the other six patients receive treatment B. The data is tabulated below:  Patient A1 Patient A2 Patient A3 Patient A4 Patient A5 Patient A6  Treatment A 0 1 0 1 1 1  Patient B1 Patient B2 Patient B3 Patient B4 Patient B5 Patient B6  Treatment B 1 0 1 0 1 0   **Comparison must be survived and by the survived by the patient B6  **Treatment B 1 0 1 0 1 0  **Comparison must be survived and survived by the survived by t	Name DEBAR	ZSHI			24	7662	1478	-
INSTRUCTIONS: Read the questions carefully and completely. Answer each question and show all your work in the space provided. Credit cannot be given if work is not shown. Good luck!  Analysis of Survival Data for Cancer Treatment (25 pointes)  A local hospital is analyzing the effectiveness of two different cancer treatments: A and B. The hospital wants to assess whether the survival rate six months after treatment is significantly different between the two treatments. Treatment is considered effective if a patient survives for at least six months after receiving it.  Data has been collected from 12 patients who received treatment, with survival after 6 months recorded as 1 (Survived) or 0 (Did not survive). 6 patients are randomly selected to receive treatment A, and the other six patients receive treatment B. The data is tabulated between treatment A and the other six patients receive treatment B. The data is tabulated between treatment A and 1 patient A2 Patient A3 Patient A4 Patient A5 Patient A6  Patient A1 Patient A2 Patient A3 Patient A4 Patient A5 Patient A6  Treatment A 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1		CHAK	RABORT	4				
A local hospital is analyzing the effectiveness of working the working the hospital wants to assess whether the survival rate six months after treatment is significantly different between the two treatments. Treatment is considered effective if a patient survives for at least six months after receiving it.  Data has been collected from 12 patients who received treatment, with survival after 6 months recorded as 1 (Survived) or 0 (Did not survive). 6 patients are randomly selected to receive treatment A, and the other six patients receive treatment B. The data is tabulated below:  Patient A1 Patient A2 Patient A3 Patient A4 Patient A5 Patient A6  Treatment A 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	all your work in the	<b>VS:</b> Read the ne space prov	questions car rided. Credit	efully and cor cannot be giv	en n work	io not bio.	question and s vn. Good luck!	show
Patient A Patient As Patient As Patient As Patient As Patient B 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	A local hospit The hospital want different between least six months a Data has been	al is analyzh s to assess w the two treat fter receiving collected fro arvived) or ( the other six	whether the suments. Treat it.  Im 12 patients (Did not suments received)	refless of two irvival rate sizement is consider who received irvive). 6 pa	x months a dered effect treatment tients are B. The da	fter treatnive if a pat , with surv randomly ta is tabula	ival after 6 mo selected to recated below:	or at onths ceive
Patient B1 Patient B2 Patient B3 Patient B4 Patient B5 Patient B6  Treatment B 1 0 1 0 1 0  L. (1pts) What are the measures the hospital took to exercise control over the study conditions?  The experiment of all patient conducted in same involved in same.  Northall same environment.  Treatments or given nandomy to patients.  All 12 participants and affected by concert.  2. (1pts) Define appropriate random variables and corresponding sample spaces.  Xi = survival status of ith patient receiving that B  Yi = survival status of ith patient receiving that B  Ci = 1(1)6)   safe space of x = 1x = 20,1)  Normival status of the data collected in this study?  Explain briefly.  We cannot me normal distribution as a model for the data collected in this study?  Normival status of the data collected since it is study?  Normival status of the data collected in this study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is study?  Normival status of the data collected since it is status of the data		Patient A1	Patient A2	Patient A3	Patient A4	Patient A		i
Treatment B 1 0 1 0 1 0  If (1pts) What are the measures the hospital took to exercise control over the study conditions?  The experiment of all patients conducted in same.  hospital same environment.  Treatments are given nandomly to patients.  All 12 participants and affected by concert.  2. (1pts) Define appropriate random variables and corresponding sample spaces.  Xi = survival status of ith patient necessary treatment A  Ci = 1(1) 6)  Yj = survival status of ith patient necessary treatment B  Yj = survival status of ith patients	Treatment A		1	0	1	1		=
Treatment B 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1 1 0 1 0 1 1 0 1 0 1 1 0 1 0 1 1 0 1 0 1 1 0 1 0 1 1 0 1 0 1 1 0 1						Patient I	0	_
Treatments are given nandomly to patients.  Treatments are given nandomly to patients.  All 12 participants and affected by concer.  2. (1pts) Define appropriate random variables and corresponding sample spaces.  Xi = survival status of ith patient necessing treatment A ci=1(1)6)  Yi = survival status of ith patient necessing that B  Yi = survival status of ith patient necessing that B  Ci=1(1)6)   safe space of x = 1x = 20,13  Cj=1(1)6)   safe space of x = 1x = 20,13  Explain briefly.  We cannot use normal distribution as a model for the data collected in this study?  Explain briefly.  We cannot use normal distribution as a binner (same model) for the data collected in this study?  Normal model for the data collected since (same model) for the data collected since (same model) for the data collected since (same model)	Treatment B	1	0	1		1		-
random variable).	A  2. (1pts) Defin  X; =  Y; =  S. (1pts) Can  Explain brie  We co	ne experience was pital ( rentmen  11 12 y e appropriate  Survivi  Ci=11  Survivi  Cj=11  we use a not fly.	ment of some some its are particles e random var val stat (1) 6) ral stat (1) 6) rmal distribu  Me The	environ  environ  given  nts are  iables and corr  is of  tion as a mod  dota	nent.  nandon  affect  responding  ith pa  ith pa  del for the  colle	and by sample spatient patient and the down of the dow	cancer.  cancer.  cancer.  receiving  traceiv  not survived  cted in this st	trustment A  ing tot B  ing tot B
<b>1</b>				1	Man	gow	VAJUMO	

For the rest of the problems, we assume that the random variables associated with the survival of each treatment are iid realizations from a Bernoulli distribution. The probability mass function (PMF) of a Bernoulli distribution is given by

$$P(Y = y; p) = p^{y}(1 - p)^{1 - y}, \quad y \in \{0, 1\}; \quad 0 (1)$$

4. (1pts) Let  $p_A$  be the probability of survival under treatment A, and  $\ell(p_A)$  be the log likelihood function of  $p_A$  based on the data from treatment A. Give a mathematical expression for  $\ell(p_A)$ .

Likelihood = LCpA) = bA I' (1-pA) n- Iyi log Likelihood =  $L(p_A) = 4 \log(p_A) + 2 \log(1-p_A)$ 5. (2pts) Let  $U(p_A)$  be the score function for  $p_A$ . Give a mathematical expression for  $U(p_A)$ .

$$U(p_A) = \frac{2}{2p_A} L(p_A)$$

$$= \frac{4}{p_A} + \frac{2}{(1-p_A)} (-1) = \frac{4}{p_A} - \frac{2}{1-p_A}$$

 $\mathcal{G}$ . (2pts) Derive the maximum likelihood estimate of  $p_A$ .

rive the maximum likelihood estimate of 
$$p_A$$
.

$$\begin{array}{ccc}
V(p_A) &= 0 &= \\
& & & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
& & \\
&$$

7 (2pts) Give a mathematical expression for the expected information in a random

size 6 from the probability density function (1).

size 6 from the probability density function (1).

$$T(p) = p^{y}(1-p)^{1-y}$$

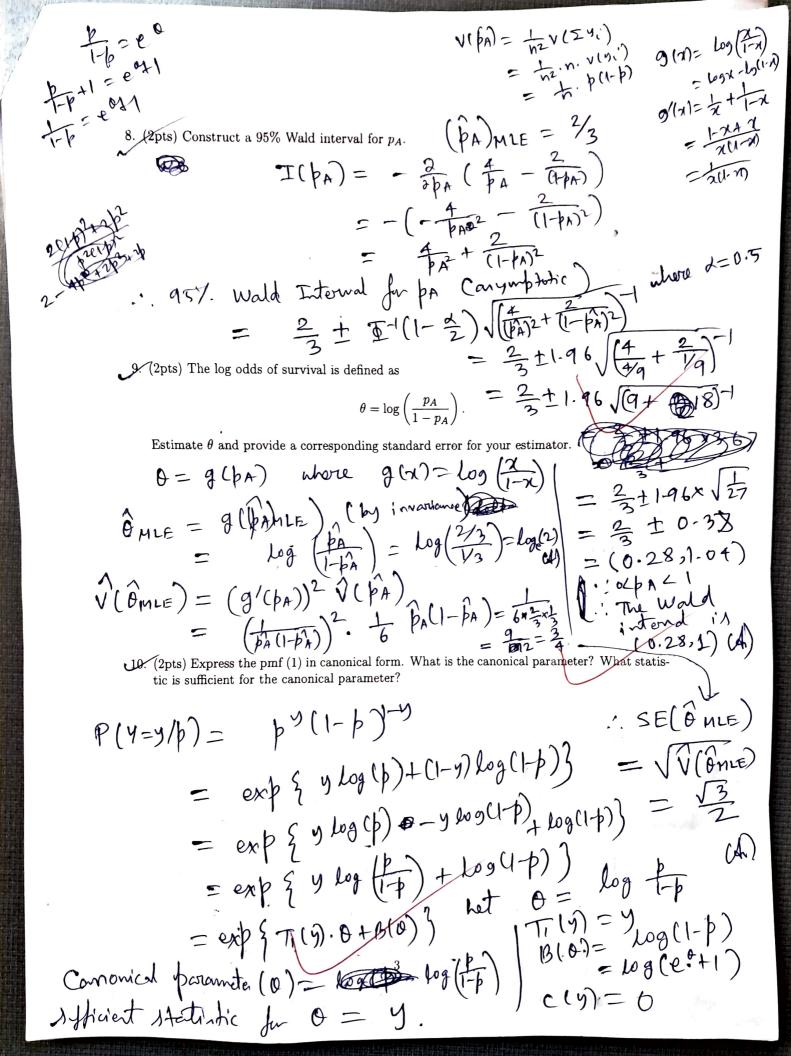
$$L(p) = p^{y}(1-p)^{1-y}$$

$$L(p) = y^{1} \log(p) + (1-y)$$

$$\log(p) = \frac{1-y}{1-p}$$

$$L(p) = \frac{1-y}{p^{2}} + \frac{1-y}{(1-p)^{2}}$$

$$L(p) = \frac{1-y}{p^{2}} + \frac{1-y}{p^{2}}$$



1. (2pts) Use your canonical parametrization to derive an expression for $E[Y]$ as a function of the canonical parameter, where Y is a Bernoulli random variable with pmf (1). Show the
steps of your derivation. We know by the properties of exponential family
$E(T(Y)) = \frac{2}{30}B(0)$
=) E(Y) = \frac{2}{20} \log(e^0+1) = \frac{e^0+1}{e^0+1} \log(-1) = \frac{1}{e^0+1}
$=\frac{y_{1-p}}{y_{1-p}}=\frac{y_{1-p}}{y_{1-p}}$
$= p \cdot CAr$

12. (1pts) Is it true that the probability density function (1) is a member of the natural exponential family? You only need to write "Yes" or "No."

13. (1pts) Is it true that the probability density function (1) is a member of the exponential dispersion family? You only need to write "Yes" or "No."

YES:
$$\begin{pmatrix}
\text{exp} & \text{for "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write "Yes" or "No."} \\
\text{vision family? You only need to write$$

14. (1pts) Is it true that the probability density function (1) forms a location/scale family? You only need to write "Yes" or "No."

## NO.

15. (2pts) Do you think that there is a significant difference in effectiveness between Treatment A and Treatment B? Justify your answer through a likelihood ratio test for the relevant parameter.

A and Treatment B? Justify your answer throughout throughout the rameter.

Ho: 
$$\beta A = \beta B$$
 $L(\beta A, \beta B) = \beta A (1-\beta A)^2 \beta B^3 (1-\beta B)^3$ 

Under  $H_1$ :  $\beta A = \frac{2}{3}$ ;  $\beta B = \frac{7}{12}$ 

Under  $H_0$ :  $\beta A = \beta B = \frac{7}{12}$ 
 $T = -2 \left( \ln(\beta) - \ln(\beta A, \beta B) \right) - \frac{1}{2} \ln(\beta A) + \frac{1}{3} \ln(\beta A)$ 
 $= -2 \left( \frac{1}{2} \ln(\beta A) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A)$ 
 $= -2 \left( \frac{1}{2} \ln(\beta A) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right) + \frac{1}{3} \ln(\beta A, \beta B)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right)$ 
 $= -2 \left( \frac{1}{3} \ln(\beta A, \beta B) + \frac{1}{3} \ln(\beta A, \beta B) \right$ 



18. (2pts) Construct a 95% Wald interval for the difference in proportions of survival between Treatment A and Treatment B. we wat would interned for by and by one ind.  $(p_A, p_B) \sim N_2(p_B), p_B$ where  $\Sigma = \left(\frac{4}{(+2)} + \frac{2}{(+2)}\right)^{-1}$ g(x,y) = x(-y) method  $\therefore \beta_A - \beta_B \longrightarrow N(\beta_A - \beta_B)(10) \sum_{i=1}^{n} (i)$ =)  $p_{A} - p_{B} \sim N(p_{A} - p_{B})(p_{A}^{2} + \frac{2}{(1-p_{A})^{2}}) + (p_{B}^{2})^{2} + (1-p_{B})^{2}$ · Wald C-I. for pA-pB is given by (pA-PB) I 1.96 (\$ 12 1 1-12) + (BB) + (BB) + (BB)  $= \left(\frac{2}{3} - \frac{1}{2}\right) \pm 1.96\sqrt{\frac{1}{23}} + \frac{1}{24} = \frac{1}{6} \pm 1.96 \times 0.28$ T=-2[7/27-7/2012+5/25+5/2512 -4/23+2/23+2/23+36/22] =-2[7/2017+5/2015-12/212+6/23+4/292]  $(1)(1-\alpha) = 5-99$   $Tobs < \chi^{2}_{cis} (1-\alpha)$  Accepted is accepted is accepted in the constant of the constanQ 6 -e-4 /2 =0.05 +115

17. (2 Extra Credit) The hospital also recorded the size of the tumor  $X_{i,j}$  for treatment i and patient j when the treatment began, and would like to investigate the relationship between the size of the tumor and the probability of survival of the patients. Write a statistical model for this investigation, and discuss what are the relevant parameter(s) to estimate, how to estimate the parameter(s), and how to make statistical inference.

Xij = rize of tumon belonging to ith fort and ith partiest j=1,2,3,4,5,6  $E(X_i) = \mu(X_{ij}) =$ E (Yi) = M2 (Xij) = p araneters 1312 are the B's.