

| Class        | Sex: Mal       | Femule            | Class | Sex Mal                    | e   Female    |
|--------------|----------------|-------------------|-------|----------------------------|---------------|
| Int          | 0              | 0                 | Int   | 118                        | 4             |
| 2nd          | 0              | 0                 |       | 154                        | 13            |
| 3rd          | 35             | 17 .              | 3rd   | 387                        |               |
| Crew         | 0              | 0                 | Crew  | 670                        | 89            |
|              | - Child, Surv  |                   |       |                            |               |
| Class        | Sex: Male      | Female            |       | = Adult, Sure<br>Sex: Male |               |
| Class<br>Int | Sex: Male      | Female 1          |       | Sex: Male                  | Female        |
| 1at<br>2nd   | Sex. Male<br>5 | Female<br>1<br>13 | Class | Sex: Male                  |               |
| Class<br>Int | Sex: Male      | Female 1          | Class | Sex: Male<br>57            | Female<br>140 |

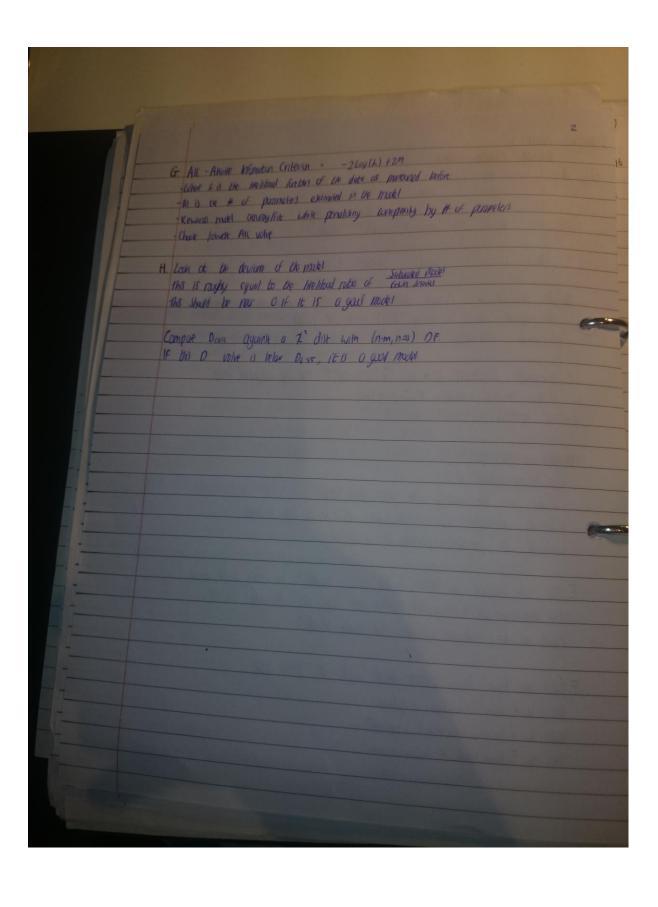
(a) Remits the information in table 1 in a new table where the first column corresponds to the group index, the response variable is in the number of people who have survived in group i, in, in the total number of people in group i, Age, in the large of group i, Sey, in the lark of group i and Citaxe, in the economic status of

Page 3 of 5

|        | ) INCOL I COMPETED ?   |
|--------|--|
| 01.21  |  |
| 7/03/1 | 6 ALSM2<br>EXAM PAPER 2014   |
| QI I   | A Exprential Distribution is a special case of the weight distribution with $\Lambda=1$  |
|        | = (1)(y)(y") exp[-04] = 0 exp[-04]   |
|        |  |
|        | B E[9] = L y P(4/1/10) dy  |
|        | = Po AVB Ay, ext [-B2,] gA   |
|        | Libstitution p = Oy du = O Ny dy   |
|        | p. vg  |
| 40     | = 60 yh u exp[-p] du   |
|        | = 200 (M/O)'M exp E-43 dN  |
|        | $= \frac{1}{6} \Gamma(1 + \frac{1}{4}) = \Gamma = \frac{1}{6} S^{N-1} \exp(-s) ds$   |
|        |  |
|        | C. 0 70 Etyl & 6" ER"  |
|        | Log: IR+ = IR for x+B  2- Sortable line fortion to map IR+ onto R  |
|        | The invert function is exp (link function is inversible)   |
|        | The inverx made is exp.  |
|        | 0-4 ER'  |
|        | -Corresponds to a duration or interpreted as duration  |
| 10-    | - Suture for multiling time  |
|        |  |
|        | The state of the s |
|        | S(T) = 1-F(t) probability of survival beyond time t  |
|        | ((1) = 1 ((c) plane)   |
|        |  |
|        | E Hozard Function  |
|        | $F(y) = 1 - exp(-\theta \tau^{\lambda})$   |
|        | - 1 1 A VIII - II I  |
|        | 11 [6(4)] 1  |
|        | H(T) = a Logistin  d (SD))  H(T) = Accelerated failure time > curit do this writing with exponential distribution  |
|        | 11/1) = Anderoted failure time? CONT 00 CM   |
|        | it is dependent on the choic of i  |

| 2 Julie 1 SERRESTER?  24 July A July 2  EARM PAPER 2014  Q2 A July A Age sow Character 1 Server 1 Serv |          |  |
|--|----------|--|
| 2 4 1 91   |          |  |
| 2 4 1 91   |          |  |
| 2 4 1 91   | , 1      | WELL I COMPLETE  |
| EAAM PAREN 2014  Q2 A 1 91   | Marie "  | MILLY I SEMESTER?  |
| EAAM PAREN 2014  Q2 A 1 91   | 29/03/16 | ALIM 2   |
| STO CHAI From Indiable   Indiab   |          | Carre Contract Contra |
| 16 20 2013 Adult Frank Grow  16 20 2013 Adult Frank Grow  16 20 2013 Adult Frank Grow  18 Ur binarrol or Power divinition to limite respect 4, 41. 41.  The respect whole is binary (Oct) with 10 graps  Binarrol Ovolution (an be used to model the farmed of each graps  C Binarrol divinition (glan as: [5] 10° (1-6)° 7 for any individual graps  P(41, 1911 101, 1610) = This p(416) 1  = This [6] 1107 (1-8) 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1   | 4- /     | 1 S StO (bill the 1st  |
| B Like binomed or Papisan distribution to Briefl religion 9, 9, 9, 9, 10  The response write is binary (Oct) with to groups  Girminal Distribution given as: [9] let (1-6) mile the burned of coin grape  C. Binomial distribution given as: [9] let (1-6) mile the burned of coin grape  P(9, 9) let e.in) = 17 in p(9,18i)  = 17 in [9] [10] mile 0 or 1  - Culd allo be sinching like yet modes forch  - Use to represent an attribute with the best levely like a battle  Indicable the population present of a wordle  E For Sex, easiest to define it as a Q1 event, O for mule, I for sinche  For any, their are 2 categories, child, which the Q1 indicable wordle  Culd use an indicable visually for levels of (10)?  F. (all we lagit half frequent of legiting) of probet firsts  MI be to best to be again a battless.  |          | 2 1 1+0 (hld Femile 10   |
| B De binned or Power dividents to make response 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4,   |          | i 1110 Child More 1"   |
| The response whole is binary (Oat) when to goups  Binamial Obstitution given as: (g) let (1-t) make the summed of each year  P(40, 900) to, 1000 = The p(10) to)  = The (g) 100 mile:  A binary variable Oat  Custo allo be sinching like year probessed.  - We be represent an attribute with the levell like a both  Indicate the presention-present of a whole feeth.  For aye, true are 2 antegories, child, which the color of indirect whole  For aye, true are 2 antegories, child, which the class.  F. Custo the light hill firsten to lay line) or probet forth  (II) be thins: the agent to steels.   |          | 16 20 20+3 Adult Fernale Crew  |
| The response whole is binary (Oat) when to goups  Binamial Obstitution given as: (g) let (1-t) make the summed of each year  P(40, 900) to, 800 = The p(1016)  = The (g) 100) 100 = The p(1016)  = The (g) 100) 100 = The p(1016)  - A binary variable Oat 1  - Custo allo be sinching like year modes force  - We be represent an attribute water the levell like a both  - Indicate the peacefron-present of a variable  E For sex, easiest to active it as a C1 event, O for mode, I for sinche  For aye, trave are 2 antegories, child, which Use Q1 indicate variable  Could use an indicator whole for levels of class.  (All the legit hill fireton of legitive) or probet forth  (III) b + bytes + by ayer + by class.  (III) b + bytes + by ayer + by class.  | R        | On brown I a do had a few to be the few to the   |
| Binomial Overham an be used to make the summed of each grap  C. Binomial distriction given as: [g] 18" (1-6)" For any individual grap  P(y, y, y, 10, y, e, 10) - The p(y, 10)  - The Giller" (1-8)" (1-8)" (1-8)"  - A binary variable  - A binary variable  - A binary variable  - Could alle be Sinething like year modef force  - Use his represent an attack water two levels like a balan  - Indicable the prometon-presence of a variable  E For sex, easiest to define it as a a, 1 events, 0 for mode, 1 for smoke  For aye, there are 2 categories, child, advist Use Q1 indicable variable  Could use an Indicable variable for levels of class  F. Could use for Indicable variable for levels of class  F. Could use for Indicable variable for levels of class  Mill by the bases + base | -        | The religinate wright is binary (Oct.) With 16 grups   |
| P(4., 96) 8., 266 = Tris p(918)  - Tris (Siller) 11 (1-8) 12-31  D Indicator Variable  - A binary variable 0 or 1  - Could allo be sirething lift gex model fearle  - Use to represent an attribute with this level lift a lization  - Indicate the presentation-present of a variable  E Foi sex, easiest to define it as a Crit event, O for male I for some  For aye, there are 2 categories, childred with Use QI indicator variable  Could use an indicator variable for levels of clust  F (allo like legit high fraction 8= lay(1) or probat fundar  MI b + byxxx + by ayer + by class:  MI b + byxxx + by ayer + by class:   |          |  |
| P(4., 96) 8., 266 = Tris p(918)  - Tris (Siller) 11 (1-8) 12-31  D Indicator Variable  - A binary variable 0 or 1  - Could allo be sirething lift gex model fearle  - Use to represent an attribute with this level lift a lization  - Indicate the presentation-present of a variable  E Foi sex, easiest to define it as a Crit event, O for male I for some  For aye, there are 2 categories, childred with Use QI indicator variable  Could use an indicator variable for levels of clust  F (allo like legit high fraction 8= lay(1) or probat fundar  MI b + byxxx + by ayer + by class:  MI b + byxxx + by ayer + by class:   |          | Pine I Indian are a [n   194 (1900) Con an inhibid area  |
| D Indicate Variable  - A binary variable 0 or 1  - Courty allow be screening like year moderflower  - Use to represent an attribute water this levell, like a both  - Indicate the premisfron-present of a variable.  E For Sex, easiest to define it as a Coll event, o for moder I for some  For aye, were are 2 actegories, child, advis Use Q1 indicate variable  (wild use an indicator variable for levels of class  F (all use layit links fraction 6= lay(100) or probat forch  Mill by t book + box ages? + boctors:  Mill by t box + box ages? + boctors:  |          | P(4  |
| - A binary variety O or 1  - Could allo be simething like yex model finish  - Use to represent an attribute with two levels like a badin  - Indicate the proposed for a variable  - Exist sex, easiest to define it as a O,1 event, O for mode, I for simule  - For aye, there are 2 categories, childradusty. Use Q1 indicate variable  - Could use an indicate variable for levels of (last)  - F. Could use fight limb function \(\theta = \left  \frac{1}{100} \) or probations.  - Finish to be six + by ayer + by alease.  - The binary variable of the levels of (last)  - The binary variable of the levels of (last)  - The binary variable of the last function of the las | -        | = 11 [9] (8) 4: (1-01) 1-4:  |
| - A binary variety 0 or 1  - Could allo be simething like yex model fonce  - Use to represent an attribute both that this levels like a bath  - Indicate the proposed for a variable.  - Exist sex, easiest to define it as a Q1 event, 0 for male, 1 for some  - For aye, there are 2 categories, childrentially Use Q1 indicate variable  - Could use an indicator variable for levels of (last):  - F. Could use facility limbs function $\theta = log(\frac{1}{lno})$ or probabilists.  - Military the age of the course.  | 1        |  |
| - Cuto allo be sinething lite yex male finale  - Use to represent an attribute with two levels, lite a body  - Indicate the present for a whole  - Indicate the present for define it as a a a consider  - For sex, easiest to define it as a a a general, o for male, I for finale  - For aye, there are 2 categories, child will be a indicator whole  - Could use an indicator whole for levels of class.  - For aye, there are 2 categories, child will be a indicator whole  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories, child will be a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For aye, there are 2 categories are a considered.  - For  |          |  |
| - Use by represent an attribute with this levell, like a below  Indicated the present for attribute with this levell, like a below  For sex, easiest to define it as a all event, o for mole, I for female  For age, there are 2 categories, child, adult Use at indicator variable  Could use an indicator variable for levels of class.  F. Could use light hink function of lay(\(\frac{1}{100}\)) or probations.  MI: by the bytes the ages \(\frac{1}{100}\) by (\(\frac{1}{100}\)).  |          | Could all be prophing life yex male/famile   |
| - Indicate the presentation-present of a variable.  E For Sex, easiest to define it as a O,1 event, O for mole I for sensule  For age, there are 2 categories, child, addity. Use Q1 indicator variable  (all we an indicator variable for levels of class.)  F. (all we light limb function $\theta$ - lay( $\frac{1}{110}$ ) or probabilists.  M1: by t bysis: + by age; 2 t by (1005);  (112: by + bysis: + by age; 2 t by (1005); 2  | -        | the to represent an attribute with this levell like a badan  |
| E For Sex, easiest to define it as a Q1 event, O for male I for senate  For aye, there are 2 categories, child, adult Use Q1 indirate variable  Could use an indivator variable for levels of clusts  F. Could use light limb fraction 0 = lay(\frac{1}{170}) or proble forch  M1: bo + by sax; + by aye; 1 by Class;  M2: bo + by sax; + by aye; 1 by Class;  | -        | ndically the present present of a varioble   |
| For age, there are 2 categories, child, adult Use 91 marcos various  Could use an indicator variable for levels of class  F. Could use legit limb function & lay(1/10) or probat function  M1: bs + bs soci + bs age; 2 + bs class;  M2: bs + bs soci + bs age; 2 + bs class;  |          |  |
| F. Call the logit link function $\theta = \log(\frac{1}{pro})$ or probatively  M1: by + by 500; + by age; 2 + by (1035);  M2: by + by 500; + by age; 2 + by (1035);  | EF       | of Sex, easiest to define it us a Oil event, O for Mole, I to know   |
| F. Call We light link frequent $\theta = \log(\frac{1}{100})$ or probationals  (M2: by + by soci + by age; 2 + by (1035);  (M2: by + by soci + by age; 2 + by (1035);  | F        | ar age, there are 2 altegories, Chilly addity the grant of dall  |
| M1: bo + bysox: + b2 (1981 + b3 (103);  M2: bo + bysox: + b2 (1981 + b2 (103)) <sup>2</sup>  | (4)      | uld we on ladicate various to levels a class   |
| M1: $b_0 + b_1 x_0 + b_2 c y_{E1} + b_3 c (c) x_1^2$ M2: $b_0 + b_1 x_0 x_1^2 + b_2 c y_{E1}^2 + b_2 c (c) x_1^2$  | - 6      | all on 1 2 half fronten A= lastin) or probat firster   |
| M2: by + by sox; + be age; 2 + by (101); 2   | r. (0    | III) UR TOUR AND TO COME A PORTOR  |
| M3 by threx: (lun: + brex: +brex;  |          |  |
|  |          | the second the base of the bas |
|  |          | D) Greenson , Garage Land  |
|  |          |  |

| 1       |  |  |  |  |  |   |        |  |
|---------|--|--|--|--|--|---|--------|--|
|         | ) lufch  | 1 SEMESTER   | 2  |  |  |   |        |  |
| 2.9103/ | /16 ALIN   |  |  |  |  |   |        |  |
| 02      | A i  | M PAYER  |  |  |  | Classi  |        |  |
|         | 7 1  | 9:   | n;<br>540  | Agei   | Sex:   | 130   |        |  |
|         | 1  | ,  | 110  | Child  | Femily.  | 1st   |        |  |
|         | 3  | li li  | 1110   | Child  | Muh  |   |        |  |
|         | 1  |  |  |  |  |   |        |  |
|         | 16   | 20   | 20+3   | Adult  | Female   | Crew  |        |  |
|         |  |  |  |  |  |   |        |  |
|         | B Ux   | bloomed or   | Passan di  | stribution b   | inulet re  | sporate 4., 42 5                              | 710    |  |
| 1       | The .  | an Louise straigh  | th IC heave 1  | (Dal) Wit  | h 10 4104  | )   |        |  |
|         | Biou   | nial Distribu  | tan can be   | und b  | model the  | survival of each                              | y y    |  |
|         |  |  |  |  |  | any individual gro                            |        |  |
|         | ( Ring   | mal distribut  | tin given as   | (4/6,11  | (b) tu   | gy morotour gro                               | 1      |  |
|         | C. Ding  |  |  | 7 12 40 30   |  |   |        |  |
|         | Ply  | , , 416) 01,   | BIG) = TI=1  | p [9: 18:1   |  |   |        |  |
|         | Ply.   | , , 416 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  | 016) = Tr=1<br>1 (8:)4: (1-8:)^  | p (4: 18:)   |  |   |        |  |
|         | P(y  | , , 416) (71, , , , ; ; ; ; ; ; ; ; ; ; ; ; ; ; ; ;  | 016) = Tr=1<br>  [8:14: (1-8:)^  | p (9: 18:1<br>:-9:   |  |   |        |  |
|         | P(y)   | , , 916 01, , = 1712 (9)   | 016) = Tr=1<br> [6:]4: [[-6:]^   | p (4:18:1<br>:-4:  |  |   |        |  |
|         | D India  | = 17 10 (G)  | 016) = Tr=1<br>  [6:)4: (1-0:)^  | p (9/16/1<br>:9:   |  |   |        |  |
|         | D India - A b  | = 1712 (2)   | 016) = 17=1<br>18:14: (1-8:)^  | 19:15:1<br>:-9:<br>uex Mulé  | /soule   |   |        |  |
|         | P(y)   | HI Variable  | O or I  on attached  | gex Mule   | /Ame   |   |        |  |
| - 20    | P(y)   | HI Variable  | O or I  on attached  | gex Mule   | /Ame   |   |        |  |
| 100     | P(y)  D India  - A bi  - Carlo  - Use I  - India                           | you to, give the Vonide inay variable to the second to the process | O or I  continue like  continue like | gex mole   | /fank<br>no levell /   | ile u kudu                                    |        |  |
| P       | P(y.  D India  - A bi  - Curlu  - Use L  - India                           | y 911 Ci.  = 1712 Ci.  uhr Vorible  inay varible  I allo be s  to represent  | O or I  On attribute  one from - presence  | gex male   | /fank<br>no levell /   | lle a lexin                                   | Clinik |  |
| 750     | P(y.  D India  - A bi  - Curlu  - Use L  - India                           | y 911 Ci.  = 1712 Ci.  uhr Vorible  inay varible  I allo be s  to represent  | O or I  On attribute  one from - presence  | gex male   | /fank<br>no levell /   | lle a lexin                                   | Clinik |  |
| 757     | P(y.  D India  - A bi  - Curlu  - Use I  - India  E For                    | you to,  = 1712 Gi  ub Vorible  inay varible  I all be s  w represent  the present  sex, easiest   | O or I  O or I  On attribute  to define its  orderiones.   | gex moles  | /finite and levelly / white I event, C UX QI   | lle a books<br>for male, I fo<br>indicate was | Clinik |  |
| 7801    | P(y.  D India  - A bi  - Curlu  - Use I  - India  E For                    | you to,  = 1712 Gi  ub Vorible  inay varible  I all be s  w represent  the present  sex, easiest   | O or I  O or I  On attribute  to define its  orderiones.   | gex moles  | /finite and levelly / white I event, C UX QI   | lle a books<br>for male, I fo<br>indicate was | Clinik |  |
| 10      | D India  A b  Call  Use I  India  E Far a  Call                            | y, un on,  - 17:5. Gi.  out Vorible:  (nay variable)  out be s  so represent  out the present  sex, easiest  y, there are  use on India  | O or I  inelling life  on attribute mellion-presente  to define it categories, co  | yex mole with the of a war of a war hild, will ultip   | /Annie white | ile u keln<br>for mote, I fo<br>indicate und  | Clinik |  |
| 750     | D India  A b  Call  Use I  India  E Far a  Call                            | y, un on,  - 17:5. Gi.  out Vorible:  (nay variable)  out be s  so represent  out the present  sex, easiest  y, there are  use on India  | O or I  inelling life  on attribute mellion-presente  to define it categories, co  | yex mole with the of a war of a war hild, will ultip   | /Annie white | ile u keln<br>for mote, I fo<br>indicate und  | Clinik |  |
| 190     | P(y.  D. India  - A b.  - Cull  - Use I  - India  E Fu  Fur a  Culd        | you to,  = 17:15 (G)  who Vonide  inay variable  I allo be s  so represent  of the process  sex, easiet  ye, there are indust  the limit him   | O or I  mething life  on attack  to define it  categories, co  | gex mole  with the cf a war  of a wa | /Annie white | ile u keln<br>for mote, I fo<br>indicate und  | Clinik |  |
| 750     | D India  A b  Culu  Use I  Indias  E For a  Culd  F. Cull 1                | He lyst lines  | Darl  Darl  Inelling like  Ch attribute  The activity  The | gex mole  with the of a war  os a a a  hild, adult  for leve  le loy live  (loss.  | /Annie white | ile u keln<br>for mote, I fo<br>indicate und  | Clinik |  |
| 100     | P(y.  D India  - A b  - Culu  - Use L  - India  E For a  Culd  For a  Culd | you to,  = 17 15 Gi  alto Vonide inay variable I allo be s  to present e to prese  sex, easielt ye, over are the layet h  to the present  the layet h  to the sex.   | D or I  melhing like  on attribute  melhing resorte  to define its  cutegories, co  cutegories | gex mole  cos a co  hildraduly  for leve  Class  22(103)2  | /Annie white | ile u keln<br>for mote, I fo<br>indicate und  | Clinik |  |
| 780     | P(y.  D India  - A b  - Culu  - Use L  - India  E For a  Culd  For a  Culd | you to,  = 17 15 Gi  alto Vonide inay variable I allo be s  to present e to prese  sex, easielt ye, over are the layet h  to the present  the layet h  to the sex.   | Darl  Darl  Inelling like  Ch attribute  The activity  The | gex mole  cos a co  hildraduly  for leve  Class  22(103)2  | /Annie white | ile u keln<br>for mote, I fo<br>indicate und  | Clinik |  |



| 4  |  |
|--|--|
| 2  | MICCH I SEMPSTALS  |
|  | The state of the s |
| 9/03/16  | ALSMI  |
| Q3 A.  | DISOIDULTION OF DOTA   |
|  | D De data a binou automo? is only worlden brothe?  |
|  | are all studying a time with survival such as in an anathroped for a day date  |
|  | the we modelling a time until failure or at we modelling the number of   |
|  | TOTAL OF CHANGE  |
| vo<br>-Fa  | edited values should all follow the retrespective distribution, and this, the edited values should all follow the retrespective distribution; any after predicted lives are not ligitually possible example, a researcher may be interested in predicting one of three possible outland this cute, the dependent variable can take only 3 distant values, and the distribution of  |
|  | dependent visions is sold to be multinarial.   |
|  | Suppok you are trying to predict people's family planning choice, specifically how y children, as a funtary of incident various other indicates  |
|  | dependent various - # children is district and most likely the distribution of that variable is  |
| high   | y snowed 111 lbn (lake, hould be reatmoste to assume dependent variable follows a poisson  |
| and the same of th |  |
| Poisson  | - # of occurates y EIN   |
| Binom  | cal - # of Surelles in nonall  |
| Expone   | trol - time until faille > survival analyti YEIR'  |
| THE RESERVE AND ADDRESS OF THE PERSON NAMED IN COLUMN  | - tine until failue YEIR !   |
| Mulaton  | rul - Binamul but with multiple (built out cure)   |
|  | The state of the s |
|  | function of the mean will be modelled as linear in the predictors?   |
| - Choire   | of linu function > must be involible   |
| - A Secu   | nd reason why the linear (multiple regression) model might not be inadequale   |
| to descr   | the a particle relationing is that the effect of the provider of the   |
| Mana   | ent roundly mus not be linear in Nature  |
|  | the positivities between a person's age or varous inaliators   |
| of ho  | ofth is most likely not linear in notice: During early adulthood   |
|  |  |
|  | TU YI  |

| 2  | J ME |
|--|------|
| the annual hollinstown of people who are to compared on that of somewer who is so  | 16.  |
|  | 10.  |
| to mornedly different between 60 to 15 probably grace.   |      |
| - the itelebrating is mon-timer in notice.   |      |
| It him between age and rather stone is the vertice   |      |
| pawer relationship in this example.  Toget by (this probet: Incorporate) or by   |      |
| Marting IR TIR Squar and Squar Pal   |      |
|  |      |
| -A smooth importable linearising link function (c) which transforms the expectation of the property condition (i) = After to a linear amplitude.   |      |
| - Chiapit was press  | An   |
| - Anades the relativity between 1/124 predicts and the men of the cliticalist fronts   |      |
| C What will the freduters go?  |      |
| - The prediction used negligible some fine training  |      |
| - As a bais, all predicts shall be inchested at ling leans allusture All.  |      |
| (reade intention) between production is in the square of a production to   |      |
| attempt to create a more accorde mules.  |      |
| - Culculde All by the parti  |      |
| - Chart Made I with lowelf All   |      |
| - Investigut periore work to determe it muld periods a war fet or not  |      |
|  |      |
| limitation of litelihud finition assumes response are independent  |      |
| COUNTY IN COUNTY OF THE O ZERO, OIL AUTH DIGITAL AGES  |      |
| - My not be rebut because of cution.   |      |
| This o remove willers in loca about  |      |
| - Same problem with hypesica apposits  | +    |
| The state of the s |      |
|  |      |
|  |      |
|  |      |
|  |      |
|  |      |
|  |      |
|  |      |