

①

F-24 DATA-240 Data Mining / Analytics
Homework - 2

Prayag Nikul Purohit (017416737)

Problem 3 Logistic regression concepts

a) eqⁿ for a logistic regression model

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

$$\begin{aligned}\log\left(\frac{p}{1-p}\right) = & -15.3001 + (0.0018 \times \text{factor 1}) \\ & + (-0.061 \times \text{factor 2}) + \\ & (0.0057 \times \text{factor 3}) + (0.0066 \times \text{factor 4}) \\ & + (0.0071 \times \text{factor 5}) + (0.1113 \times \text{factor 6}) \\ & + (-0.0098 \times \text{factor 7}) + (0.0686 \times \text{age})\end{aligned}$$

b) ① $Z = \frac{\text{coefficient}}{\text{std error}}$

intercept $\Rightarrow Z = \frac{-15.3001}{0.0015} = \boxed{-10,200.067}$

factor 1 $\Rightarrow Z = \frac{0.0018}{0.0103} = \boxed{0.17476}$

$$\text{factor 2} \Rightarrow z = \frac{-0.0061}{0.0105} = \boxed{-0.5809}$$

$$\text{factor 3} \Rightarrow z = \frac{0.0057}{0.0028} = \boxed{2.0357}$$

$$\text{factor 4} \Rightarrow z = \frac{0.0066}{0.0038} = \boxed{1.7368}$$

$$\text{factor 5} \Rightarrow z = \frac{0.0071}{0.0199} = \boxed{0.3568}$$

$$\text{factor 6} \Rightarrow z = \frac{0.1113}{0.0492} = \boxed{2.2622}$$

$$\text{factor 7} \Rightarrow z = \frac{-0.0098}{0.0037} = \boxed{-2.6456}$$

$$\text{Age} \Rightarrow z = \frac{0.0686}{0.0224} = \boxed{3.0625}$$

$$\textcircled{ii} \quad p = 2 \times (1 - \text{CDF}(|z|))$$

Interupt

$$p = 2 \times (1 - \text{CDF}(1020.0667))$$

$$\boxed{p = 0}$$

z is extremely large so p -value is essentially 0

factor 1

$$p = 2 \times (1 - \text{CDF}(0.1748))$$
$$= \boxed{0.86}$$

(2)

factor 2

$$p = 2 \times (1 - \text{CDF}(0.5810))$$
$$= \boxed{0.56}$$

factor 3

$$p = 2 \times (1 - \text{CDF}(2.0357))$$
$$= \boxed{0.0418}$$

factor 4

$$p = 2 \times (1 - \text{CDF}(1.7368))$$
$$= \boxed{0.0824}$$

factor 5

$$p = 2 \times (1 - \text{CDF}(0.3568))$$
$$= \boxed{0.7213}$$

factor 6

$$p = 2 \times (1 - \text{CDF}(2.2622))$$
$$= \boxed{0.0237}$$

factor 7

$$p = 2 \times (1 - \text{CDF}(2.6486))$$
$$= \boxed{0.0081}$$

Age

$$p = 2 \times (1 - \text{CDF}(2.30125))$$
$$= \boxed{0.0022}$$

(iii) Lower bound = $\beta - 1.96 \times SE$
 upper bound = $\beta + 1.96 \times SE$ } for 95% CI

$$CI = [LB, UB]$$

Intercept

$$LB = -15.3001 - (1.96 \times 0.015)$$

$$= -15.30304$$

$$UB = -15.3001 + (1.96 \times 0.015)$$

$$= -15.29716$$

$$CI = [-15.30304, -15.29716]$$

factor 1

$$LB = 0.0018 - (1.96 \times 0.0103) = -0.018388$$

$$UB = 0.0018 + (1.96 \times 0.0103) = 0.021988$$

$$CI = [-0.018388, 0.021988]$$

factor 2

$$LB = -0.0061 - (1.96 \times 0.0105) = -0.0266$$

$$UB = -0.0061 + (1.96 \times 0.0105) = 0.01448$$

$$CI = [-0.0266, 0.01448]$$

factor 3

$$LB = 0.0057 - (1.96 \times 0.0028) = 0.000212$$

$$UB = 0.0057 + (1.96 \times 0.0028) = 0.011188$$

$$CI = [0.000212, 0.011188]$$

factor 4

$$LB = 0.0066 - (1.96 \times 0.0038) = -0.000848$$

$$UB = 0.0066 + (1.96 \times 0.0038) = 0.014048$$

$$CI = [-0.000848, 0.014048]$$

factor 5

$$LB = 0.0071 - (1.96 \times 0.0199) = -0.031904$$

(3)

$$UB = 0.0071 + (1.96 \times 0.0199) = 0.046104$$

$$CI = [-0.0319, 0.0461]$$

factor 6

$$LB = -0.0098 - (1.96 \times 0.0492) = -0.01486$$

$$UB = -0.0098 + (1.96 \times 0.0492) = 0.20773$$

$$CI = [0.01486, 0.20773]$$

factor 7

$$LB = -0.0098 - (1.96 \times 0.0037) = -0.01705$$

$$UB = -0.0098 + (1.96 \times 0.0037) = -0.002548$$

$$CI = [-0.01705, -0.00254]$$

Age

$$LB = 0.0686 - (1.96 \times 0.0224) = 0.024696$$

$$UB = 0.0686 + (1.96 \times 0.0224) = 0.112504$$

$$CI = [0.024696, 0.112504]$$

cy intercept

$$p\text{-value} = 0.00 \text{ (significant)}$$

$$CI = [-15.303, -15.297] \Rightarrow \text{no '0' in it}$$

$$\Rightarrow \text{Significance} \rightarrow \text{Yes}$$

factor 1

$$p\text{-value} = 0.861 \text{ (no)}$$

$$CI = [-0.018, 0.0219] \Rightarrow \text{(include '0')}$$

$$\Rightarrow \text{Significance} \rightarrow \text{No}$$

factor 2 { $p\text{-value} = 0.5613$ (not significant)
 $CI = [-0.0266, 0.01448] \rightarrow$ (includes '0')
Significance \rightarrow No

factor 3 { $p\text{-value} = 0.0418$ (significant)
 $CI = [0.0012, 0.01118] \rightarrow$ (no '0' in it)
Significance \rightarrow Yes

factor 4 { $p\text{-value} = 0.0824$ — (not significant)
 $CI = [-0.00848, 0.014] \rightarrow$ (includes '0')
Significance \rightarrow No

factor 5 { $p\text{-value} = 0.721$ (not significant)
 $CI = [0.0148, 0.2077] \rightarrow$ (no '0' in it)
Significance \rightarrow Yes

factor 6 { $p\text{-value} = 0.0237$ (significant)
 $CI = [0.0148, 0.2077] \rightarrow$ (no '0' in it)
Significance = Yes

Age
factor 7 { $p\text{-value} = 0.0022$ (significant)
 $CI = [0.02469, 0.1125] \rightarrow$ (no '0' in it)
Significance = Yes

factor 7 $\left\{ \begin{array}{l} p\text{-value} = 0.0081 - (\text{significant}) \\ CI = [-0.017, -0.0025] - (\text{no '0' in it}) \end{array} \right.$ (4)
Significant - Yes

Analysis of each factor

1) $p\text{-value} < 0.05$ } for 95% of confidence
2) CI including '0' }

2) Age \Rightarrow $\text{coeff} = 0.0686 \Rightarrow +ve$ so if age \uparrow
the log-odds of events (death due to the disease) \uparrow
 \Rightarrow 1 unit increase of age, odds increases by 0.0686

Also it is a significant feature as $p\text{-value} < 2$
& CI doesn't have '0' in it.

Factor 7 \Rightarrow $\text{coeff} = -0.0098 \Rightarrow -ve$ if factor 7 \uparrow
log-odds decreases
 \Rightarrow 1 unit increase of factor 7, the log-odds decreases by 0.0098.

Also it is a significant feature as
 $p\text{-value} < 2$ & CI don't have '0' in it

7(i) a 100 unit decrease in factor 7

$$\beta = -0.0098 \times 100 = -0.98$$

$$e^{-0.98} \approx \underline{0.3753}$$

$$100 \times (1 - 0.3753) = \underline{62.47\%}$$

a 100 unit decrease in factor 7 is associated with a 62.47% reduction in odds of death.

ii) an additional year of age

$$\beta = 0.0686 \times 1 = 0.0686$$

$$e^{0.0686} = 1.071$$

$$(1.071 - 1) \times 100 = \underline{7.1\%}$$

an year of age associated with 7.1% increase in odds of death.

$$b) \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n$$

(5)

$$= -15.3001 + (0.0018 \times 125) + (-0.0061 \times 105) \\ + (0.0057 \times 235) + (0.0066 \times 105) + \\ (0.0071 \times 12.5) + (0.1113 \times 42) + \\ (-0.0098 \times 475) + (0.0686 \times 50)$$

$$= -9.81085$$

$$p = \frac{1}{1 + e^{-(-9.81085)}} = \frac{1}{1 + 18298.748} \\ = \frac{1}{18299.748}$$

$$p = 0.0005485$$

\Rightarrow So probability of death for 50 year old woman with the given factor values is 0.00548%