

Generalized Linear Models Homework

Unity_id - ssamudr

Question1)

After training the logistic regression model for all predictors using the training data, we can infer that “Category_Coins/Stamps” is the predictor with the highest estimate for its coefficient. We then train model.single using only “Category_Coins/Stamps” as the predictor. The equation of relation of response variable “Competitive” to “Category_Coins/Stamps” is:

$$Competitive = -0.507 - 2.07233 * CategoryCoins/Stamps$$

- a) The relationship of log-odds with the predicted variable is:

$$\log (P(Y = yes)/(1 - P(Y = yes))) = B_o + B_1 X_1$$

$$P(Y) = 1/(1 + e^{B_0 + B_1 x})$$

$$P(Y = yes) = 1/\left(1 + e^{(-0.5 - 2.072 \cdot categorycoins/stamps)}\right)$$

- b) $\log (P(Y = yes)/(1 - P(Y = yes))) = B_o + B_1 X_1$
 $\log (odds) = -0.507 - 2.072 \cdot categorycoins/stamps$
 $odds = e^{-0.5 - 2.072 \cdot categorycoins/stamps}$

- c) $\text{logit} = \log (odds) = -0.5 - 2.072 \cdot categorycoins/stamps$

Question 2)

The top four predictors for model_all model are: Category_Coins/Stamps, Category_Health/Beauty, Category_Clothing/Accessories, currency_GBP. Based on these predictors, we can express the different equations as follows:

a)

$$\text{Logit} = \log(\text{odds}) = -0.5 - 2.072 * \text{Category_Coins/Stamps} - 1.64 * \text{Category_Health/Beauty} - 1.42 * \text{Category_Clothing/Accessories} + 1.26 * \text{currency_GBP}.$$

b)

$$\text{Odds} = e^{(-0.5 - 2.072 * \text{Category_Coins/Stamps} - 1.64 * \text{Category_Health/Beauty} - 1.42 * \text{Category_Clothing/Accessories} + 1.26 * \text{currency_GBP})}$$

c)

$$P(Y=\text{yes}/X) = 1/(1+e^{\text{logit}}) = 1/(1+e^{(-0.5 - 2.072 * \text{Category_Coins/Stamps} - 1.64 * \text{Category_Health/Beauty} - 1.42 * \text{Category_Clothing/Accessories} + 1.26 * \text{currency_GBP})})$$

Question 3)

The highest predictor for model_all is Category_Coins/Stamps. The generalised equation for odds can be given as:

$$\text{Odds} = e^{(-0.5 - 2.072 * \text{Category_Coins/Stamps} - 1.64 * \text{Category_Health/Beauty} - 1.42 * \text{Category_Clothing/Accessories} + 1.26 * \text{currency_GBP})}$$

Now, if we increase the value of Category_Coins/Stamps by 1 and keep the value of coefficients constant, the corresponding equation for odds is:

$$\text{Odds}' = e^{(-0.5 - 2.072 * (\text{Category_Coins/Stamps} + 1) - 1.64 * \text{Category_Health/Beauty} - 1.42 * \text{Category_Clothing/Accessories} + 1.26 * \text{currency_GBP})}$$

Hence the odds ratio is:

$$\text{Odds/Odds}' = e^{(-0.5 - 2.072 * \text{Category_Coins/Stamps} - 1.64 * \text{Category_Health/Beauty} - 1.42 * \text{Category_Clothing/Accessories} + 1.26 * \text{currency_GBP})} / e^{(-0.5 - 2.072 * (\text{Category_Coins/Stamps} + 1) - 1.64 * \text{Category_Health/Beauty} - 1.42 * \text{Category_Clothing/Accessories} + 1.26 * \text{currency_GBP})}$$

$$\text{odds}'/\text{odds} = e^{-2.072}$$

Therefore, if the value of Category_Coins/Stamps increases by 1, the value of the response variable changes by a factor of $e^{-2.072}$. If it was linear regression, then the value of response would change by the factor of 2.072 (coefficient) times. Since, the value of logistic regression gives us the output of the logit function, we have to find the value of e . Whereas in linear regression the value coefficient output directly affects the response.

Question 4 :

We can use anova test to check if the two models- fit_reduced and fit_all are equivalent to each other or not. After running anova test on the two models. The p_value obtained is 0.7086 . Since p-value is greater than 0.05, we can conclude that the two models do not significantly differ from each other.

Question 5 :

Overdispersion occurs when the value of $\phi = \gg 1$. In this case, Residual deviance Residual df $\phi = 0.992$. Since the value of ϕ is close to 1, there is no overdispersion present in the model. If there is any overdispersion, then we can use quasi-binomial distribution instead of binomial family distribution.