

TOPIC 2: Generalized Linear Models

1. X_h : predictor with the highest estimated (absolute value) of regression coeff.

(a) Probabilities: $\text{Prob}(Y = \text{Yes} | X_h = X)$

$$\rightarrow E(Y) = \text{logit}(P) = \log_e \left| \frac{P}{1-P} \right| = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_q X_q$$

$$P = \frac{1}{1 + e^{-Z}}$$

$$\text{where } Z = \beta_0 + \beta_1 X.$$

In our case,

$X = \text{Category everything else}$: a dummy variable
or

$$X = \text{Category}$$

So,

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 (\text{Category}^*))}}$$

$$(b) \text{ Odds : } \frac{M_y}{1-M_y} = \frac{P}{1-P} = e^{\beta_0 + \beta_1 (\text{Category}^*)}$$

$$(c) \text{ Logit : } \log_e \left(\frac{P}{1-P} \right) = \beta_0 + \beta_1 (\text{Category}^*)$$

* Here we can also use $\text{Category everything else}$ as well. But since it is a dummy category

2 fit all model with the top 4 predictors in terms of absolute value estimates
(2) The four predictors with the highest absolute value estimates are:

X_1 : Category \rightarrow Category everything else

X_2 : currency \rightarrow currency GBP

X_3 : endDay ~~Sat-Wed~~ \rightarrow endDay Sat-Wed

X_4 : openforce

(i) Probability, $P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4)}}$

(ii) Odds: $e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4}$

(iii) Logit: $\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$

3. X_n : predictor with highest estimate in fit.all
 $X_n = \text{category everything else}$ or Category

Since its a categorical predictor it can take only 2 values
 0 or 1.

So assuming $X_n = 1$ originally, we change to 0 by
 doing a single unit increase.

$$\frac{\text{odds}(X_n = 1)}{\text{odds}(X_n = 0)} = \frac{e^{\beta_0 + \beta_1 X_n + \dots + \beta_q X_q}}{e^{\beta_0 + \beta_1 X_n + \dots + \beta_q X_q}}$$

$$\Rightarrow \frac{\cancel{e^{\beta_0}} \cdot e^{\beta_1 X_n} \cdot \cancel{e^{\beta_2 X_2}} \dots \cancel{e^{\beta_q X_q}}}{\cancel{e^{\beta_0}} \cdot e^{\beta_1 X_n} \cdot \cancel{e^{\beta_2 X_2}} \dots \cancel{e^{\beta_q X_q}}}$$

$$\Rightarrow \frac{e^{\beta_1 X_1}}{e^{\beta_1 X_0}} \Rightarrow \frac{e^{\beta_1}}{1} \Rightarrow e^{\beta_1}$$

$$\beta_1 \text{ for category everything else} = -2.041$$

$$\text{So odds ratio} = e^{-2.041} = 0.1298$$

4. We built a reduced model with 6 significant predictors (everything except "Duration").

We then carried the anova(), chi-square version of the test.

The test gave us a non-significant chi-square value of 0.2783. which suggests that fit.reduced fits as well as fit.all.

5. We can test for overdispersion by computing the mean deviance.

$$\text{Mean Deviance} = \frac{\text{Residual Deviance}}{\text{Residual Degree of Freedom}} = \frac{1189.6}{1167} = 1.019$$

Since mean deviance is close to 1, the model is well fitted.

THE END