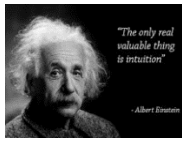


Multinomial Logistic Models

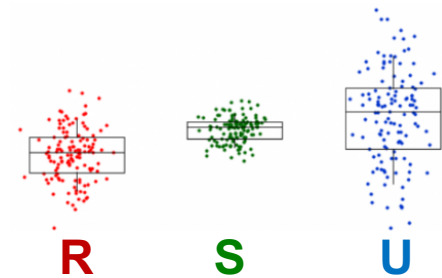


KOGOD SCHOOL *of* BUSINESS
AMERICAN UNIVERSITY • WASHINGTON, DC



Multinomial Logistic: Intuition

- In **binomial logit** the response variable has **2 possible values** (e.g., “good mail” and “spam”), so we can easily quantify this variable with 2 dummy variables coded **0** (e.g., good mail) and **1** (e.g., spam)
- If the response variable has **more than 2 categories** we have a **multinomial logistic** situation
- For example, if we are trying to find the type of **home location** a person is likely to purchase (e.g., **rural**, **suburban**, **urban**).
- We can solve this problem **fitting 2 binomial models**.
- We first select the “**reference**” category (e.g., rural) that will be coded **0** in **both models**.
- We then fit **two binomial models**, one for **suburban vs. rural** and another for **urban vs. rural**.
- More generally, a response variable with **K categories** can be modeled with **K-1 binomial** logistic models. The category left out is called the “**reference**” category.
- **Multinomial logistic** regression **fits** all the binomials **together**



Multinomial Logistic Details

- A **multinomial** model can be solved by **fitting binomial logistic models one by one**, but this yields separate fit statistics for model
- A **multinomial** logistic model **fits** all the binomials **together** and provides fit statistics for the model as a whole
- The response variable Y has K categories
 - We select one of them as the “**reference**” category k_R
 - We estimate logistic **coefficients** for the other $K-1$ categories
 - In **binomial** logistic, a coefficient for X_I represents the increase in log odds of $Y=1$, relative to $Y=0$, when X_I increases by 1 unit.
 - In **multinomial** logistic, the coefficient for a predictor X_{Ik} for the k^{th} category represents the increase in log odds of $Y=k$, relative $Y=k_R$, when X_I increases by 1 unit.

Multinomial Logistic Interpretation

For a **multinomial logit** model:

$$\begin{aligned}
 \text{Logit}(Y) = & \beta_{01} + \beta_{11}(X_1) + \beta_{21}(X_2) + \beta_{31}(X_3) + \dots \\
 & \beta_{02} + \beta_{12}(X_1) + \beta_{22}(X_2) + \beta_{32}(X_3) + \dots \\
 & \dots \\
 & \beta_{0K} + \beta_{1K}(X_1) + \beta_{2K}(X_2) + \beta_{3K}(X_3) + \dots + \varepsilon
 \end{aligned}$$

- Y has K categories (e.g., Freshman, Sophomore, Junior, Senior withdrawing from school)
- $\text{Logit}(Y)$ is the **log odds** of the response variable being in category k (e.g., Junior), **relative** to the **reference** category k_R (e.g., Freshman)
- β_{ik} is the **effect** of variable X_i (e.g., GPA) on the log odds of $Y = k$
- But **odds** and log odds are always **relative** to something (e.g., leaving vs. staying)
- So, β_{ik} measures how much the log odds Y of being k , relative to the reference category k_R (e.g., Freshman), **change** when X_i (e.g., GPA) goes up by 1 unit.



Multinomial Logit: Fit Statistics

- The **vglm(){VGAM}** function in **R** reports the:
 - **Log-Likelihood**
 - **Deviance** (2LL) $\rightarrow -2 * \text{Log Likelihood}$
 - **AIC** \rightarrow calculated by adding $2 * \text{Number of Variables}$

Confusion Matrix:

$$\text{Error Rate} = \frac{\text{Incorrect}}{\text{Total}} = \frac{\text{Off - Diagonal}}{T}$$

$$\text{Sensitivity}_{\text{class}} = \frac{\text{Correct}_{\text{class}}}{\text{Total}_{\text{class}}} = \frac{TP_A}{A_A}; \frac{TP_B}{A_B}; \frac{TP_C}{A_C}$$

$$\text{Specificity}_{\text{class}} = \frac{\text{Correct}_{\text{NotClass}}}{\text{Total}_{\text{NotClass}}} =$$

$$\frac{TP_B + TPC}{A_B + AC} \text{ (for A)}; \frac{TP_A + TPC}{A_A + AC} \text{ (for B)}; \frac{TP_A + TPB}{A_A + AB} \text{ (for C)}$$

Pred	Actual			Total
	A	B	C	
A	TP _A	P _{A/B}	P _{A/C}	P _A
B	P _{B/A}	TP _B	P _{B/C}	P _B
C	P _{C/A}	P _{C/B}	TP _C	P _C
Total	A _A	A _B	A _C	T

Tips

`vglm()` {VGAM} → “Vector Generalized Linear Model” function in the {VGAM} “Vector Generalized and Additive Model” package is a function to fit **multinomial logistic**; “vector” refers to the fact that the response variable is no longer binary, but a vector of categorical values.

Note: there are other packages that can fit multinomial logit models, with their own strengths, difficulties and challenge for you to explore, e.g.: `multinom()` {nnet} and `glmnet()` {glmnet}

`vglm.fit =`
`vglm(y~x1+x2+etc., family=multinomial(refLevel=1),`
`data=dataName)` → `refLevel=1` indicates the first category in the response variable will be left out as a baseline

Like binary logistic with `glm()` → `logLik(vglm.fit); -`
`2*logLik(vglm.fit); deviance(vglm.fit);` and `AIC(vglm.fit)`
 work well with a `vglm()` object



KOGOD SCHOOL
of
BUSINESS

