# Postestimation Analysis with Stata by SPost13 commands of Survey Data analyzed by MNLM

Debora Giovannelli – debora.giovannelli@gmail.com

## 1-Introduction

Data coming back from a brand survey have been analyzed by a regression model for nominal outcomes, also known as the Multinomial Logit Model.

The Multinomial Logit Model (MNLM) belongs to a multivariate version of Generalized Linear Models (GLM), a class of models popularized by McCullagh and Nelder (1982) and widely used in many different fields (Social Sciences, Biomedical Sciences, Epidemiology, Public Health, Genetic, Zoology, Education, but also Marketing Researches, Survey Analysis and Product/Process/Service Quality Control).

The interpretation of these models requires a background knowledge that is not always common, especially in business application fields.

**Data must be "readable" to anyone who has the responsibility to take serious decision**, which can strongly influence not only the business of a company but also the safety and the quality of its products/processes and services.

The scope of this presentation is to show and highlight the advantages of the implementation of SPost13 commands, setup by J. Scott Long and J. Freese, as very useful tools for making easier the interpretation of results coming from the implementation of this regression model for nominal response variables.

## 2-Objectives

The interpretation of regression models for categorical response variables is complex because of their nonlinearity.

Models for nominal and ordinal outcomes may be interpreted using odds ratios (for logit models) and quantities based on predicted probabilities (*predictive margins*). While odds ratios do not depend on the values of the predictors (multiplicative effects), the meaning of odds ratios in terms of probabilities depends on the values of all the regressors (so the magnitude of probability change depends on the values of all the explicative variables in the model).

Because of nonlinearity these models require postestimation analysis and computation of predicted probabilities and related quantities as marginal effects, in order to fully describe the effects of all predictors.

**Methods** for the **interpretation** of *nonlinear regression models* for *categorical outcomes* have been proposed by *J. Scott Long* and *Jeremy Freese* [7].

The statistical analyses here referred have been implemented by *Stata®/15.1* and *SPost13* (Stata postestimation commands for version 13), a suite of programs for the postestimation interpretation of regression models for categorical outcomes, developed by J.S. Long and J. Freese, with the object to **give evidence on how SPost13 postestimation commands make easier the interpretation of nonlinear models as the MNLM**.

## 3-Dataset Description and Explorative Data Analysis

These statistical analyses are based on data coming from a **survey** conducted for assessing **Customer orientation** in the **professional audio market**, and previously analyzed by modelling the probability of respondents choice (favourite brand selection) by a Multinomial Logit Model (**MNLM**), where some characteristics of the respondents where included as explicative variables [6]. The response variable **Brand** is a multilevel nominal categorical variable with **5 outcomes** (5 brands coded A, B, C, D, Others), while the two categorical explanatory variables, specified in the model, are a binary variable *X1* (*Age*), with two levels (age over 50 years old, age under or equal to 50 years old), and a categorical variable *X2* (*PVM, Primary Vertical Market*) with 4 levels: Ent (*Entertainment*), GER (*Government Institution, Educational, Religious Institutions*), Oth (*others*), R&S (*Rental & Staging*).

*Variables description*
. codebook Brand X1 X2

**Brand**
```
type:   numeric (long)
label:  Brand
range:  [1,5]                units:  1
unique values:  5           missing .:  0/741
tabulation:  Freq.  Numeric  Label
              243      1      A
              156      2      B
               45      3      C
              194      4      D
              103      5      Others
```

**X1**
```
type:   numeric (long)
label:  X1
range:  [1,2]                units:  1
unique values:  2           missing .:  0/741
tabulation:  Freq.  Numeric  Label
              398      1      Over 50
              343      2      ≤50
```

**X2**
```
type:   numeric (long)
label:  X2
range:  [1,4]                units:  1
unique values:  4           missing .:  0/741
tabulation:  Freq.  Numeric  Label
              163      1      Ent
              249      2      GER
              161      3      Oth
              168      4      R&S
```

*Three-way cross-tabulation table*
. table X2 Brand, by(X1) center stubwidth(12)

| Age and PVM | A | B | C | D | Others |
|---|---|---|---|---|---|
| **Over 50** | | | | | |
| Ent | 32 | 15 | 3 | 31 | 8 |
| GER | 55 | 13 | 4 | 35 | 23 |
| Oth | 35 | 8 | 5 | 27 | 22 |
| R&S | 26 | 17 | 8 | 20 | 11 |
| **≤50** | | | | | |
| Ent | 17 | 28 | 5 | 14 | 10 |
| GER | 24 | 31 | 8 | 37 | 11 |
| Oth | 24 | 11 | 3 | 14 | 11 |
| R&S | 22 | 33 | 9 | 16 | 6 |

## 4-Model fitting and Selection

*Estimation using mlogit command*

The MNLM has been fit using **mlogit** command.

The dependent variable Brand has 5 nominal outcomes (A, B, C, D, Others).

The model has been parameterized setting **category A** as **base outcome** (reference group).

The independent variables, both categorical, have been included in the model by using **factor-variable notation**.

Four models have been fitted:
- Full model **"mfull"**: with two regressors with interaction terms (saturated model)
- Main model **"mmain"**: model with two regressors X1 (Age) and X2 (PVM) but no interaction terms
- Restricted model **"mX1"**: model with the regressor X1 (X2 omitted)
- Restricted model **"mX2"**: model with the regressor X2 (X1 omitted)

The following table summarizes the Information Criteria for all fitted models.
. estimates stats m*

Akaike's information criterion and Bayesian information criterion

| Model | Obs | ll(null) | ll(model) | df | AIC | BIC |
|---|---|---|---|---|---|---|
| mfull | 741 | -1103.296 | -1063.711 | 32 | 2191.421 | 2338.877 |
| mmain | 741 | -1103.296 | -1069.208 | 20 | 2178.415 | 2270.555 |
| mX1 | 741 | -1103.296 | -1085.362 | 8 | 2186.724 | 2223.588 |
| mX2 | 741 | -1103.296 | -1085.931 | 16 | 2203.862 | 2277.59 |

*Estimation results for the main effects model*
. mlogit Brand ib(1).X1 ib(2).X2, base(1) vsquish nolog

```
Multinomial logistic regression          Number of obs  =    741
                                          LR chi2(16)    =
                                          Prob > chi2    = 0.0000
Log likelihood = -1069.2077               Pseudo R2      = 0.0309
```

| Brand | Coef. | Std. Err. | z | P>|z| | [95% Conf. Interval] |
|---|---|---|---|---|---|---|
| **A** | (base outcome) | | | | | |
| **B** | | | | | | |
| X1 ≤50 | 1.095761 | .2167432 | 5.06 | 0.000 | .6709519 | 1.52057 |
| X2 Ent | .604048 | .2848636 | 2.12 | 0.034 | .0457256 | 1.16237 |
| GER | -.3798146 | .3271769 | -1.16 | 0.246 | -1.02107 | .2614403 |
| Oth | -.7379691 | .2794202 | -2.61 | 0.009 | -1.801156 | 1.275623 |
| R&S | -1.286163 | .228511 | -5.63 | 0.000 | -1.734036 | -.8382894 |
| **C** | | | | | | |
| X1 ≤50 | .6500734 | .3296256 | 1.97 | 0.049 | .0040191 | 1.296128 |
| X2 Ent | .2005366 | .4915536 | 0.41 | 0.683 | -.7628908 | 1.163964 |
| GER | -.0258795 | .4882302 | 0.05 | 0.958 | -.9310341 | .9827931 |
| Oth | .9463166 | .4199749 | 2.26 | 0.024 | .124976 | 1.767657 |
| R&S | -2.304349 | .3590409 | -6.42 | 0.000 | -3.008056 | -1.600642 |
| **D** | | | | | | |
| X1 ≤50 | .1081333 | .1966082 | 0.55 | 0.582 | -.2772116 | .4934782 |
| X2 Ent | .1092642 | .2610328 | 0.42 | 0.676 | -.4023311 | .6208595 |
| GER | -.1677803 | .2586654 | -0.65 | 0.517 | -.6747552 | .3391945 |
| Oth | -.0978976 | .2720755 | -0.36 | 0.719 | -.6311558 | .4353607 |
| R&S | -.2357516 | .1805266 | -1.31 | 0.192 | -.5895773 | .118074 |
| **Others** | | | | | | |
| X1 ≤50 | -.0348331 | .2428998 | -0.14 | 0.886 | -.5109081 | .4412418 |
| X2 Ent | -.0635457 | .3420584 | -0.19 | 0.853 | -.7339678 | .6068765 |
| GER | .3861687 | .2958074 | 1.31 | 0.192 | -.1936031 | .9659405 |
| Oth | -.0986129 | .3472282 | -0.28 | 0.776 | -.7791676 | .5819419 |
| R&S | -.9251729 | .2255382 | -4.10 | 0.000 | -1.36722 | -.4831262 |

## 5-SPost13 command fitstat

*Postestimation SPost13 command fitstat*

The main effect model has been compared versus the full model by the SPost13 command fitstat.

```
. quietly mlogit Brand ib(1).X1 ib(2).X2 b(1).X1#b(2).X2, base(1)
. quietly fitstat, save
. quietly mlogit Brand ib(1).X1 ib(2).X2, base(1)
. fitstat, diff
```

| | Current | Saved | Difference |
|---|---|---|---|
| **Log-likelihood** | | | |
| Model | -1069.208 | -1063.711 | -5.497 |
| Intercept-only | -1103.296 | -1103.296 | 0.000 |
| **Chi-square** | | | |
| D(df=721/709/12) | 2138.415 | 2127.421 | 10.994 |
| LR(df=16/28/-12) | 68.176 | 79.170 | -10.994 |
| p-value | 0.000 | 0.000 | 0.529 |
| **R2** | | | |
| McFadden | 0.031 | 0.036 | -0.005 |
| McFadden(adjusted) | 0.013 | 0.007 | 0.006 |
| Cox-Snell/ML | 0.088 | 0.101 | -0.013 |
| Cragg-Uhler/Nagelkerke | 0.093 | 0.107 | -0.014 |
| Count | 0.358 | 0.364 | -0.007 |
| Count(adjusted) | 0.044 | 0.054 | -0.010 |
| **IC** | | | |
| AIC | 2178.415 | 2191.421 | -13.006 |
| AIC divided by N | 2.940 | 2.957 | -0.018 |
| BIC(df=20/32/-12) | 2270.555 | 2338.877 | -68.302 |

Note: Likelihood-ratio test assumes current model nested in saved model.

Difference of 68.302 in BIC provides very strong support for current model.

SPost13 command *fitstat* summarizes in a single table many fit statistics for comparing competing models.

## 6-SPost13 command listcoef

*Comparisons across categories by listcoef*
. listcoef, pvalue(0.05) positive

```
mlogit (N=741): Factor change in the odds of Brand (P<0.05)
```

**Variable: 2.X1 (sd=0.499)**

| | | b | z | P>|z| | e^b | e^bStdX |
|---|---|---|---|---|---|---|
| B | vs A | 1.0958 | 5.056 | 0.000 | 2.991 | 1.728 |
| B | vs D | 0.9876 | 4.384 | 0.000 | 2.685 | 1.637 |
| B | vs Others | 1.1306 | 4.213 | 0.000 | 3.097 | 1.758 |
| C | vs A | 0.6501 | 1.972 | 0.049 | 1.916 | 1.383 |

**Variable: 1.X2 (sd=0.415)**

| | | b | z | P>|z| | e^b | e^bStdX |
|---|---|---|---|---|---|---|
| B | | 0.6040 | 2.120 | 0.034 | 1.830 | 1.285 |

**Variable: 3.X2 (sd=0.413)**

| | | b | z | P>|z| | e^b | e^bStdX |
|---|---|---|---|---|---|---|
| Others | vs B | 0.7660 | 2.065 | 0.039 | 2.151 | 1.372 |

**Variable: 4.X2 (sd=0.419)**

| | | b | z | P>|z| | e^b | e^bStdX |
|---|---|---|---|---|---|---|
| B | vs A | 0.7280 | 2.605 | 0.009 | 2.071 | 1.357 |
| B | vs D | 0.8259 | 2.802 | 0.005 | 2.284 | 1.413 |
| B | vs Others | 0.8266 | 2.257 | 0.024 | 2.285 | 1.414 |
| C | vs A | 0.9463 | 2.258 | 0.024 | 2.576 | 1.487 |
| C | vs D | 1.0442 | 2.430 | 0.015 | 2.841 | 1.549 |
| C | vs Others | 1.0449 | 2.171 | 0.030 | 2.843 | 1.549 |

SPost13 command *listcoef* provides in a single table the estimates for all the comparisons of outcome categories for each variable included in the model. By specific options the output may be suitably simplified.

## 7-SPost13 command mlogitplot



SPost13 command *mlogitplot* provides a plot that synthetizes the effects of all regressors on all contrasts in odds ratio scale and in logit scale, giving also evidence of their significance

## 8-Interpretation in terms of Odds Ratios

Based on this graph, we may conclude the following:
- for individuals with age ≤50, compared with individuals with age over 50, holding PVM constant, the odds of selecting brand C or B relative to brand A significantly increase by a factor of 1,92 for C and 2,99 for B, while for the other contrasts (D vs A and Others vs A) the effects are not significant
- for individuals with PVM = Ent, compared with individuals with PVM = GER, holding Age constant, just one contrast (B vs A), is statistically significant, with an increase by a factor of 1,83
- for individuals with PVM = Oth, compared with individuals with PVM = GER, holding Age constant, all the contrasts respect to A category, are statistically not significant
- for individuals with PVM = R&S, compared with individuals with PVM = GER, holding Age constant, the odds of selecting brand C and brand B relative to brand A significantly increase by a factor of 2,58 for C and 2,07 for B, while the odds of selecting the brands D and Others relative to brand A do not significantly change

Moreover this graph provides evidence on the effects for all the other contrasts (different base outcomes).

As an example:
- for individuals with PVM = R&S, when compared with individuals with PVM = GER, holding Age constant, the odds of selecting brand C, rather than one of the other brands, is significant for the contrasts C vs A, C vs D and C vs Others, while the contrast C vs B is not statistically significant (as provided by listcoef command output).

## 9-Interpretation based on Adjusted Predictions and Marginal Effects

The **MNLM** is linear in the logit but is **nonlinear in probability**: while the factor change in the odds is constant across the levels of all variables, the effect of each predictor on the probability of an outcome of the response depends on the value (for continuous predictors) or level (for categorical predictors) of the specific predictor and on the level of all the other independent variables specified in the model (marginal effects depends on the values of all variables). This makes the interpretation complex so that the evaluation of the effects of all the explicative variables for all the logits, just based on the estimated coefficients, represents a limitation.

Moreover, models for nominal outcomes are even more complex because they provide more parameters to interpret respect to the models for ordinal outcomes, where constraints are imposed (the effect of each regressor is constrained to be equal in all equations).

The interpretation using predictions as Predictive Margins or Adjusted Predictions and summary measures based on predictions as Marginal Effects is more informative for assessing the impact of each independent variable on each outcome of the response variable.

*Adjusted Predictions* and *Marginal Effects* computation is divided in two main types: **predictions** and *summary* **margins**.

Margins provides the different types of Marginal Effects (three different approaches of computation), which depends on the different ways of controlling for the other variables in the model while computing Adjusted Predictions:
- Average Marginal Effects (**AMEs**) are computed as difference between two Average Adjusted Predictions (**AAPs**)
- Marginal Effects at Means (**MEMs**) are computed as difference between two Adjusted Predictions at Means (**APMs**)
- Marginal Effects at Representative values (**MERs**) are computed as difference between two Adjusted Predictions at specific values of the other variables (**APRs**)

In this specific context where all regressors specified in the model are categorical, the use of **factor-variable notation** in model specification is critical in order to **guarantee correct results by using margins** (this way Stata recognizes any interdependencies between variables).

Moreover, considering the categorical nature of both regressors, two types of marginal effects, *AMEs* and *MERs*, have been **computed as statistics to interpret the effect of the characteristics of the respondents on the choice of the favourite brand**.

## 10-SPost13 command mtable for tabulating Predictive Margins

*Table of Adjusted Predictions*

SPost13 command mtable is a wrapper of margins: it uses margins for building tables of Adjusted Predictions and tables of Marginal Effects (dydx).

If the outcome has multiple categories, mtable automatically submits multiple margins commands for all outcomes and combines the results in the table. Results from multiple calls of mtable may be combined into a single compact table:

```
. quietly mtable, at(X2 = 1 X1 = 2 ) rowm(PVM Ent Age ≤50) dec(4) below
. quietly mtable, at(X2 = 2 X1 = 2 ) rowm(PVM GER Age ≤50) dec(4) below
. quietly mtable, at(X2 = 3 X1 = 2 ) rowm(PVM Oth Age ≤50) dec(4) below
. quietly mtable, at(X2 = 4 X1 = 2 ) rowm(PVM R&S Age ≤50) dec(4) below
. quietly mtable, at(X2 = 1 X1 = 1 ) rowm(PVM Ent Age Over 50) dec(4) below
. quietly mtable, at(X2 = 2 X1 = 1 ) rowm(PVM Ger Age Over 50) dec(4) below
. quietly mtable, at(X2 = 3 X1 = 1 ) rowm(PVM Oth Age Over 50) dec(4) below
. mtable, at(X2 = 4 X1 = 1 ) rowm(PVM R&S Age Over 50) dec(4) below
```

Expression: Pr(Brand), predict(outcome())

| | A | B | C | D | Others |
|---|---|---|---|---|---|
| PVM Ent Age ≤50 | 0.2447 | 0.3700 | 0.0572 | 0.2402 | 0.0879 |
| PVM GER Age ≤50 | 0.3048 | 0.2519 | 0.0583 | 0.2683 | 0.1167 |
| PVM Oth Age ≤50 | 0.3258 | 0.1842 | 0.0634 | 0.2425 | 0.1836 |
| PVM R&S Age ≤50 | 0.2299 | 0.3936 | 0.1133 | 0.1833 | 0.0798 |
| PVM Ent Age Over 50 | 0.3471 | 0.1755 | 0.0423 | 0.3059 | 0.1292 |
| PVM Ger Age Over 50 | 0.3902 | 0.1078 | 0.0390 | 0.3083 | 0.1547 |
| PVM Oth Age Over 50 | 0.3933 | 0.0743 | 0.0403 | 0.2627 | 0.2294 |
| PVM R&S Age Over 50 | 0.3442 | 0.1970 | 0.0885 | 0.2466 | 0.1237 |

SPost13 command *mtable* allows tabulating Adjusted Predictions for multiple outcomes in a compact table. By multiple calls of mtable and suitable labelling options, a synthetic and informative table may be provided.

## 11-SPost13 command mchange for AMEs computation

**AMEs** are marginal effects computed as difference between two Average Adjusted Predictions (**AAPs**).

. mchange

```
mlogit: Changes in Pr(y) | Number of obs = 741
Expression: Pr(Brand), predict(outcome())
```
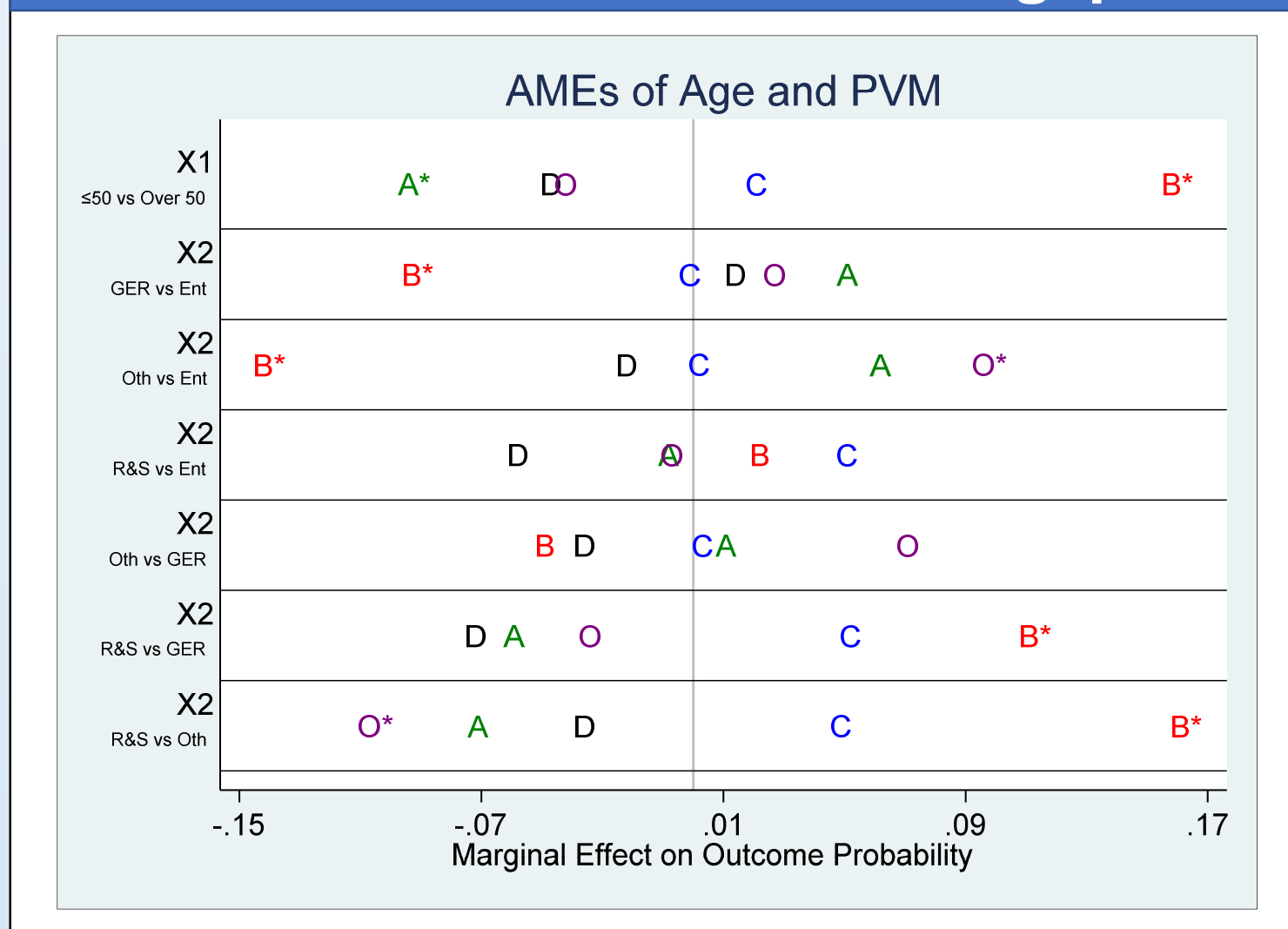
| | A | B | C | D | Others |
|---|---|---|---|---|---|
| **X1** | | | | | |
| ≤50 vs Over 50 | -0.092 | 0.160 | 0.021 | -0.047 | -0.042 |
| p-value | 0.007 | 0.000 | 0.248 | 0.149 | 0.097 |
| **X2** | | | | | |
| GER vs Ent | 0.051 | -0.091 | -0.001 | 0.014 | 0.027 |
| p-value | 0.275 | 0.027 | 0.952 | 0.752 | 0.409 |
| Oth vs Ent | 0.063 | -0.140 | 0.002 | -0.022 | 0.098 |
| p-value | 0.230 | 0.001 | 0.934 | 0.650 | 0.014 |
| R&S vs Ent | -0.008 | 0.022 | 0.051 | -0.058 | -0.007 |
| p-value | 0.867 | 0.639 | 0.076 | 0.219 | 0.843 |
| Oth vs GER | 0.011 | -0.049 | 0.003 | -0.036 | 0.071 |
| p-value | 0.814 | 0.163 | 0.881 | 0.416 | 0.066 |
| R&S vs GER | -0.059 | 0.113 | 0.052 | -0.072 | -0.034 |
| p-value | 0.199 | 0.006 | 0.063 | 0.093 | 0.294 |
| R&S vs Oth | -0.071 | 0.163 | 0.049 | -0.036 | -0.105 |
| p-value | 0.170 | 0.000 | 0.093 | 0.443 | 0.008 |

Average predictions

| | A | B | C | D | Others |
|---|---|---|---|---|---|
| Pr(y|base) | 0.328 | 0.211 | 0.061 | 0.262 | 0.139 |

## 12-SPost13 command mchangeplot



SPost13 command *mchangeplot* provides a plot that synthetizes in terms of AMEs the effects of all regressors on all contrasts, giving also evidence of their significance

## 13-SPost13 command mtable for MERs computation

**MERs** are marginal effects computed as difference between two Adjusted Predictions for a variable, conditioning at specific values of the other variables (**APRs**).

These conditional MEs may be computed by mchange or mtable SPost13 commands.

When computing MERs of Age by mchange, conditioning on PVM, multiple calls of mchange are required, because just one value of PVM can be specified in at() option.

To synthetize all the MERs in one single table **multiple calls of mtable have been submitted**.

When computing MERs for PVM by mtable, conditioning on Age, two tables have been provided in order not to loose the labels of table's rows (one for the estimates and one for the p-values).

With multiple outcomes multiple calls of *SPost13 command mtable* may provide more synthetic output of MERs

*MERs of PVM*

MEs of PVM at specific levels of Age

Expression: Marginal effect of Pr(Brand), predict(outcome())

| | A | B | C | D | Others |
|---|---|---|---|---|---|
| Ent | -0.0431 | 0.0677 | 0.0034 | -0.0024 | -0.0256 |
| Oth | 0.0030 | -0.0335 | 0.0013 | -0.0456 | 0.0747 |
| R&S | -0.0460 | 0.0892 | 0.0496 | -0.0617 | -0.0310 |
| Ent | -0.0601 | 0.1181 | -0.0011 | -0.0281 | -0.0288 |
| Oth | -0.0677 | 0.0057 | -0.0258 | 0.0468 | |
| R&S | -0.0749 | 0.1416 | 0.0550 | -0.0848 | -0.0369 |

Specified values of covariates

| | X1 |
|---|---|
| | 1 |
| | 2 |
| Set 1 Current | 1 2 |

p-values for MEs of PVM at specific levels of Age

Expression: Marginal effect of Pr(Brand), predict(outcome())

| | A | B | C | D | Others |
|---|---|---|---|---|---|
| Ent | 0.3952 | 0.0305 | 0.8542 | 0.9605 | 0.4933 |
| Oth | 0.9522 | 0.1617 | 0.9407 | 0.3262 | 0.0754 |
| R&S | 0.3645 | 0.0063 | 0.0448 | 0.1842 | 0.4015 |
| Ent | 0.1634 | 0.0273 | 0.9653 | 0.5121 | 0.3023 |
| Oth | 0.6518 | 0.1674 | 0.8356 | 0.5540 | 0.0618 |
| R&S | 0.0709 | 0.0773 | 0.0345 | 0.1744 | |

Specified values of covariates

| | X1 |
|---|---|
| Set 1 Current | 1 2 |

*MERs of Age*

MEs of Age at specific levels of PVM

Expression: Marginal effect of Pr(Brand), predict(outcome())

| | A | B | C | D | Others |
|---|---|---|---|---|---|
| PVM Ent d Pr(y) | -0.1025 | 0.1945 | 0.0148 | -0.0657 | -0.0412 |
| PVM Ent p | 0.0028 | 0.0000 | 0.3380 | 0.0579 | 0.0624 |
| PVM Ger d Pr(y) | -0.0854 | 0.1441 | 0.0193 | -0.0400 | -0.0380 |
| PVM Ger p | 0.0194 | 0.0000 | 0.1960 | 0.2561 | 0.1421 |
| PVM Oth d Pr(y) | -0.0675 | 0.1099 | 0.0237 | -0.0202 | -0.0459 |
| PVM Oth p | 0.0754 | 0.0001 | 0.1698 | 0.5397 | 0.2003 |
| PVM R&S d Pr(y) | -0.1143 | 0.1966 | 0.0247 | -0.0631 | -0.0439 |
| PVM R&S p | 0.0009 | 0.0000 | 0.3834 | 0.0360 | 0.0426 |

## 14-Interpretation in terms of Marginal Effects

For the dichotomous variable Age, which generates 1 dummy, one AME and four MERs are provided.

For PVM, which has four categories and generates three dummies, three AMEs and 6 MERs are provided.

With AMEs we are comparing two hypothetical populations. This means that the AME for the dichotomous independent variable Age compares two populations, measuring the average of the differences in adjusted predictions for two groups, one made of individuals all over 50 years old, the other made of individuals all less or equal to 50 years old, that have the same values of PVM (holding PVM as observed). Because the only difference between the two populations is their age, the changes in probability for all the outcomes may be attributed to the impact of the variable Age.

In terms of AMEs we may conclude the following:
- **on average, being ≤50**, compared with being >50, holding PVM as observed, **significantly increases the probability of selecting B** as favourite brand by **0,16** (p < 0.001) and **significantly decreases the probability of selecting A** as favourite brand by **0,092** (p < 0.01), while for the other brands the change in probability is not significant
- for the variable PVM we observe significant AMEs just for some contrasts relative to outcomes A and B. As an example, for individuals with **PVM = R&S**, compared with individuals with PVM = GER, holding Age as observed, **the probability of selecting B as favourite brand significantly increases** by **0,113** (p < 0.01)

AMEs have the limitation to provide just one single estimate for each contrast (one for Age and three for PVM).

In order to give better evidence on how the probability for each outcome may vary with the characteristics of the respondents, MERs have been computed by multiple calls of mtable to better assess the variability in effects across cases.

In terms of MERs we may conclude the following:
- the table "MERs of Age" refers the MEs of Age for different levels of PVM, showing that, when changing **from older to younger individuals**, the changes in probability to select the brand as favourite, **all significantly increases** for brand **B** and **all significantly decreases** for brand **A**
- the two tables provided for the estimates and relative *p-values* of "MERs of PVM" show **significant MEs mostly for brand B**

The **output of mtable is limited to the contrasts provided by the specified model**: to obtain **all the possible contrasts** multiple calls of mchange may be run.

## 15-Conclusions

Methods of interpretation using marginal effects for nonlinear models are provided by the Stata command margins, which allows to compute Adjusted Predictions and Marginal Effects (AMEs, MEMs, MERs).

In this poster SPost13 commands provided by J. Scott Long and Jeremy Freese have been **run by Stata** to **make easier the interpretation of a MNLM** implemented to analyze data coming back from a brand survey.

## 16-References

1. Agresti A. 2015. *Foundations of Linear and Generalized Linear Models*. John Wiley & Sons, Inc.
2. Agresti A. 2013. *Categorical Data Analysis*. 3rd ed. John Wiley & Sons, Inc.
3. Agresti A. 2018. *An Introduction to Categorical Data Analysis*, 3rd ed., John Wiley & Sons, Inc.
4. Agresti A. 2018. *Statistical Methods for the Social Sciences*, 5th edition, Pearson.
5. Jann B. 2013. *Predictive Margins and Marginal Effects in Stata*. University of Bern, 11th German Stata Users Group meeting. University of Potsdam.
6. Giovannelli, D. 2017. *Approccio Statistico all'analisi dei dati di ritorno di una Brand Survey condotta in ambito Audio Professionale*. Quality & Engineering, vol. 1(2-3): 153-186.
7. Long, J.S., and J. Freese. 2014. *Regression Models for Categorical Dependent Variables Using Stata*. 3rd ed. College Station, TX: Stata Press.
8. Long, J.S. 2014. *New methods of interpretation using marginal effects for nonlinear models*. EUSMEX 2016: Mexican Stata Users Group meeting.
9. Rising B. 2013. *Using Predictive Margins to Make Clearer Explanations*. Indian Stata Users Group meeting.
10. Williams, R. 2012. *Using the margins command to estimate and interpret adjusted predictions and marginal effects*. The Stata Journal, 12(2): 308-331.
11. Williams, R. 2019. *Using the spost13 commands for adjusted predictions and marginal effects with binary dependent variables*. University of Notre Dame, https://www3.nd.edu/~rwilliam/ Last revised January 29, 2019.
12. Williams, R. 2019. *Adjusted Predictions & Marginal Effects for Multiple Outcome Models & Commands (including ologit, mlogit, oglm, & gologit2)*. University of Notre Dame, https://www3.nd.edu/~rwilliam/ Last revised January 29, 2019.
13. Williams, R. 2019. *Multinomial Logit Models – Overview*. University of Notre Dame, https://www3.nd.edu/~rwilliam/ Last revised February 7, 2019.
14. Williams, R. 2019. *Post-Estimation Commands for MLogit*. University of Notre Dame, https://www3.nd.edu/~rwilliam/ Last revised February 7, 2019.