

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/307654108>

# An Application on Multinomial Logistic Regression Model

Article in Pakistan Journal of Statistics and Operation Research · March 2012

DOI: 10.1234/pjsor.v8i2.234

CITATIONS

36

READS

2,661

1 author:



[Abdalla M. El-Habil](#)

Al-Azhar University - Gaza

32 PUBLICATIONS 116 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



name of paper is An Application on Multinomial Logistic Regression Model [View project](#)

# An Application on Multinomial Logistic Regression Model

Abdalla M. EL-HABIL

Head of the Department of Applied Statistics

Faculty of Economics and Administrative Sciences

Al-Azhar University, Gaza - Palestine

abdalla20022002@yahoo.com

## Abstract

This study aims to identify an application of Multinomial Logistic Regression model which is one of the important methods for categorical data analysis. This model deals with one nominal/ordinal response variable that has more than two categories, whether nominal or ordinal variable. This model has been applied in data analysis in many areas, for example health, social, behavioral, and educational. To identify the model by practical way, we used real data on physical violence against children, from a survey of Youth 2003 which was conducted by Palestinian Central Bureau of Statistics (PCBS). Segment of the population of children in the age group (10-14 years) for residents in Gaza governorate, size of 66,935 had been selected, and the response variable consisted of four categories. Eighteen of explanatory variables were used for building the primary multinomial logistic regression model. Model had been tested through a set of statistical tests to ensure its appropriateness for the data. Also the model had been tested by selecting randomly of two observations of the data used to predict the position of each observation in any classified group it can be, by knowing the values of the explanatory variables used. We concluded by using the multinomial logistic regression model that we can able to define accurately the relationship between the group of explanatory variables and the response variable, identify the effect of each of the variables, and we can predict the classification of any individual case.

**Keywords:** Multinomial logistic regression model - categorical data analysis - maximum likelihood method - generalized linear models -classification.

## I. Introduction

In recent years, specialized statistical methods for analyzed categorical data have increased, particularly for application in biomedical and social science. Regression analysis is one of these statistical tools that utilize the relationship between two or more variables. The regression models can be divided into two groups, the first related to linear relationship models, and the second related to non-linear relationship models. The linear models, considered up to this point, are satisfactory for most regression applications. Nonlinear model used when the linear model is not suitable anyhow. Many of statisticians believe that the logistic regression model is one of the important models can be applied to analyze a categorical data; this model is a special case of generalized linear models (GLM). The multinomial logistic regression (MLR) model used in generally effective where the response variable is composed of more than two levels or categories. The basic concept was generalized from binary logistic regression. Continuous variables are not used as response variable in logistic regression, and only one response variable can be used. The MLR model can be used to predict a response variable on the basis of continuous and/or categorical explanatory variables to determine the percent of variance in the response variable explained by the explanatory variables, to rank the relative importance of independents, to

assess interaction effects, and to understand the impact of covariate control variables. The MLR model allows the simultaneous comparison of more than one contrast, that is, the log odds of three or more contrasts are estimated simultaneously, Garson (2009). The logistic regression model assumes that the categorical response variable has only two values, in general, 1 for success and 0 for failure. The logistic regression model can be extended to situations where the response variable has more than two values, and there is no natural ordering of the categories. Natural ordering can be treated as nominal scale, such data can be analyzed by slightly modified methods used in dichotomous outcomes, and this method is called the multinomial logistic. The impact of predictor variables is usually explained in terms of odds ratios. Logistic regression applies maximum likelihood estimation after transforming the dependent into a logit variable (the natural log of the odds of the dependent occurring or not). Logistic regression calculates changes in the log odds of the dependent, not changes in the dependent itself as ordinary least square (OLS) regression does. Logistic regression has many analogies to OLS regression: logit coefficients correspond to b coefficients in the logistic regression equation, the standardized logit coefficients correspond to beta weights, and a pseudo R square ( $R^2$ ) statistic is available to summarize the strength of the relationship. Unlike OLS regression, however, logistic regression does not assume linearity of relationship between the independent variables and the dependent, does not require normally distributed variables, does not assume homoscedasticity, and in general has less stringent requirements. It does, however, require that observations be independent and that the independent variables be linearly related to the logit of the dependent. The predictive success of the logistic regression can be assessed by looking at the classification table, showing correct and incorrect classifications of the dichotomous, ordinal, or polytomous dependent. Goodness-of-fit tests such as the likelihood ratio test are available as indicators of model appropriateness, as is the Wald statistic to test the significance of individual independent variables.

The idea of this study focusing on MLR model, that we believe it is important and useful for analyzing categorical data. Therefore, the problem is:

By using real data, how can we apply a new statistical method (multinomial logistic regression model) for analyzing categorical data?

The paper is organized as follows: Section 2 recalls the technical background of multinomial logistic regression model. Section 3 physical violence data. Section 4 building of multinomial logistic model. Section 5 conclusion.

## **II. Multinomial logistic regression model**

### **The logit (logistic) regression model**

In fact, the multinomial logistic regression (MLR) model is a fairly straightforward generalization of the binary model, and both models depend mainly on logit analysis or logistic regression. Logit analysis in many ways is the natural

complement of ordinary linear regression whenever the response is categorical variable. When such discrete variables occur among the explanatory variables they are dealt with by the introduction of one or several (0, 1) dummy variables, but when the response variable belongs to this type, the regression model breaks down. Logit analysis provides a ready alternative.

For a response variable  $Y$  with two measurement levels (dichotomous) and explanatory variable  $X$ , let:  $\pi(x) = p(Y = 1 | X = x) = 1 - p(Y = 0 | X = x)$ , the logistic regression model has linear form for logit of this probability

$$\text{Logit}[\pi(x)] = \log\left(\frac{\pi(x)}{1 - \pi(x)}\right) = \alpha + \beta x, \quad \text{where the odds} = \frac{\pi(x)}{1 - \pi(x)},$$

The odds =  $\exp(\alpha + \beta x)$ , and the logarithm of the odds is called logit, so

$$\text{Logit}[\pi(x)] = \log\left(\frac{\pi(x)}{1 - \pi(x)}\right) = \log[\exp(\alpha + \beta x)] = \alpha + \beta x$$

The logit has linear approximation relationship, and logit = logarithm of the odds. The parameter  $\beta$  is determined by the rate of increase or decrease of the S-shaped curve of  $\pi(x)$ . The sign of  $\beta$  indicates whether curve ascends ( $\beta > 0$ ) or descends ( $\beta < 0$ ), and the rate of change increases as  $|\beta|$  increases.

### **Multiple logistic regressions**

The logistic regression can be extending to models with multiple explanatory variables. Let  $k$  denotes number of predictors for a binary response  $Y$  by

$x_1, x_2, \dots, x_k$ , the model for log odds is

$$\text{Logit}[P(Y = 1)] = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

And the alternative formula, directly specifying  $\pi(x)$ , is

$$\pi(x) = \frac{\exp(\alpha + \beta_1 x_1 + \dots + \beta_k x_k)}{1 + \exp(\alpha + \beta_1 x_1 + \dots + \beta_k x_k)}$$

The parameter  $\beta_i$  refers to the effect of  $x_i$  on the log odds that  $Y = 1$ , controlling other  $x_j$ , for instance,  $\exp(\beta_i)$  is the multiplicative effect on the odds of a one-unit increase in  $x_i$ , at fixed levels of other  $x_j$ .

If we have  $n$  independent observations with  $p$ -explanatory variables, and the qualitative response variable has  $k$  categories, to construct the logits in the multinomial case, one of the categories must be considered the base level and all the logits are constructed relative to it. Any category can be taken as the base level, so we will take category  $k$  as the base level. Since there is no ordering, it is apparent that any category may be labeled  $k$ . Let  $\pi_j$  denote the multinomial

probability of an observation falling in the  $j^{\text{th}}$  category, to find the relationship between this probability and the  $p$  explanatory variables,  $X_1, X_2, \dots, X_p$ , the multiple logistic regression model then is

$$\log \left[ \frac{\pi_j(x_i)}{\pi_k(x_i)} \right] = \alpha_{0i} + \beta_{1j}x_{1i} + \beta_{2j}x_{2i} + \dots + \beta_{pj}x_{pi},$$

Where  $j = 1, 2, \dots, (k-1)$ ,  $i = 1, 2, \dots, n$ . Since all the  $\pi$ 's add to unity, this reduces to

$$\log(\pi_j(x_i)) = \frac{\exp(\alpha_{0i} + \beta_{1j}x_{1i} + \beta_{2j}x_{2i} + \dots + \beta_{pj}x_{pi})}{1 + \sum_{j=1}^{k-1} \exp(\alpha_{0i} + \beta_{1j}x_{1i} + \beta_{2j}x_{2i} + \dots + \beta_{pj}x_{pi})},$$

For  $j = 1, 2, \dots, (k-1)$ , the model parameters are estimated by the method of ML. Practically, we use statistical software to do this fitting, Chatterjee and Hadi (2006).

### Baseline-Category Logit Model

In MLR model, the estimate for the parameter can be identified compared to a baseline category. We defined bold letter as matrix or vector, let

$\pi_j(\mathbf{x}) = p(Y = j | \mathbf{x})$  at a fixed setting  $\mathbf{x}$  for explanatory variables, with  $\sum_j \pi_j(\mathbf{x}) = 1$ ,

for observations at that setting, we treat the counts at the  $J$  categories of  $Y$  as multinomial with probabilities,  $\{\pi_1(\mathbf{x}), \dots, \pi_J(\mathbf{x})\}$ , logit models pair each response category with a baseline category, often the most common model is:

$$\log \frac{\pi_j(\mathbf{x})}{\pi_J(\mathbf{x})} = \alpha_j + \boldsymbol{\beta}'_j \mathbf{x},$$

where  $j = 1, \dots, (J-1)$ , simultaneously describes the effects of  $\mathbf{x}$  on these  $(J-1)$  logits, the effects vary according to the response paired with the baseline, these  $(J-1)$  equations determine parameters for logits with other pairs of response

categories. Since  $\log \frac{\pi_a(\mathbf{x})}{\pi_b(\mathbf{x})} = \log \frac{\pi_a(\mathbf{x})}{\pi_J(\mathbf{x})} - \log \frac{\pi_b(\mathbf{x})}{\pi_J(\mathbf{x})}$

with categorical predictors, Pearson chi-square statistic  $\chi^2$  and the likelihood ratio chi-square statistic  $G^2$  goodness-of-fit statistics provide a model check when data are not sparse. When an explanatory variable is continuous or the data are sparse, such statistics are still valid for comparing nested models differing by relatively few terms, Agresti (2002).

### III. Physical violence data

We used real data of Youth Survey, 2003 which conducting by Palestinian Center Bureau of Statistics (PCBS) for application of the MLR model. The data

were used for the purposes of scientific research according to a special agreement between PCBS and Al – Azhar University- Gaza. In this survey, youth were defined as the young people, their age was (10-24) years, includes preteen, teen-agers (10-19) and young (20-24). Also in this survey age with complete years were used, the data referred to calculating the vital rates and ages as it was on 1/3/2003, the target population of this survey was consist of all Palestinian households that usually reside in the Palestine Territory.

### **Sample size and design frame used on the Youth Survey**

The sample size had been 5,570 households of whom 4,830 households responded, 3,256 households residing in the West Bank, 1,574 households in Gaza Strip. The sample strata had been designed on two levels, first level was the governorate (16 governorates), and second level was the type of locality (urban, rural and camps). The survey frame was a list of enumeration areas peculiar to the 1997 Population, Housing and Establishment Census, enumeration areas was a residential area containing about 150 households. It should be noted that the data had been collected according to the procedures, rules and methodology established by PCBS to achieve the highest data-quality, PCBS user guide (2003).

Our resource of data was the file "Youthfile.sav", conducting by PCBS of the Youth Survey, 2003, and the user guide, survey questionnaire, and methodology book. According to this survey, total number of persons in Palestinian National Authority on the age group (10-24 Years) in the year 2003 was (1,189,282), 51.0% male, 49.0% female, and 62.3% in West Bank, and 37.7% in Gaza Strip.

### **Target population group used**

The target group used in this study was youth people living in Gaza governorate in the age group 10-14 years in both sex (male and female). To achieve this goal we selected population group with the following specifications: first we selected the youth group in Gaza governorate (160,948), and then, we selected age group (10-14 Years), (66,935 or 41.588 % of the total).

### **Selection of the variables**

#### **Response variable**

Through the review of a questionnaire of the survey, there were two questions drew our attention directly related to the issue of physical violence, first question was "have you been subjected to physical violence (beating, burning, biting, pushing, etc) during the last month? ", with two levels of measures (yes/no), second was " who did practice physical violence against you?", with 10 levels of measures, or 10 kinds of people exercised of physical violence against the youth: father, mother, sibling, wife, other relatives, teacher, employer, peer (schoolmate, neighbor, etc), Israeli forces, and others.

In fact, the aim of this analysis does not focus basically on the phenomenon of violence, which it is very important, and its need special study, but the aim is to apply our statistical model, MLR model, on categorical data. For this, we tried to choose available related data on physical violence from Youth Survey 2003 according to our criteria already mentioned in the sample size and frame. By merging the two variables we got the response variable with 11 levels of measures, as the target population was the children 10-14 years, living in Gaza, wife and Israeli forces were excluded as there is no frequency to these levels. We tried to focus on the practice of physical violence on children by the family, father, mother, sibling, and very close environment of the child, schoolmates, the neighbors, and others. We note that a small number of these levels, in the same time due to skewness in the response variable we combined these levels: other relatives, teacher, employer, and other, to be in one category, Takagi et. al. (2007).

We note that, the skewness before merging was (2.908) and, and after merging (1.887), and standard error of the skewness was 0.009 for both. The response variable became as the following "have you been subjected to physical violence during the last month, and who did practice physical violence against you?", had four categories, 0-had not been, 1-father/ mother, 2-sibling, 3-peer & other, we called the response variable as "response physical violence by". The frequencies of response variable according to these categories are shown in Table 3.1.

**Table 3.1: The frequencies of the response variable categories**

| Physical violence categories | Frequency | Percent |
|------------------------------|-----------|---------|
| 0 Had not been               | 52545     | 78.5    |
| 1 Father/Mother              | 4201      | 6.3     |
| 2 Sibling                    | 7025      | 10.5    |
| 3 Peer & other               | 3163      | 4.7     |
| Total                        | 66935     | 100.0   |

#### **Baseline category (reference) of the response variable**

Any category of response variable can be chosen to be the baseline or reference category, the model will fit equally well, achieving the same likelihood and producing the same fitted values, only the values and interpretation of the parameters will change, Schafer (2006). In our situation we used the category with highest frequency so we selected category of (0-had not been). This means, the comparison will be against the children whom did not been suffering of the physical violence in the last month of survey date.

#### **The explanatory variables**

We tried to select a set of explanatory variables, that, we believed it has an effect in somehow, on the physical violence against children in age 10-14 years in Gaza in the year 2003. Some of these explanatory variables describe the environment around the child, some belongs to child himself. We will review these explanatory variables in some details.

A set of questions talk about "In your opinion, are the following behavior exists among youth in the locality where lives?":

- *Hh04-a "Alcohol consumption"* with two categories (1-no/ little, 2-yes widely).
- *Hh04-b "Smoking"* with two categories (1-no/ little, 2-yes widely).
- *Hh04-c "Reckless driving"* with two categories (1-no/ little, 2-yes widely).
- *Hh04-d "Drug abuse"* with two categories (1-no/ little, 2-yes widely).
- *Hh04-e "Verbal violence" (e.g., harassment, swearing)*, with two categories (1-no/ little, 2-yes widely).
- *Hh04-f "Begging"* with two categories (1-no/ little, 2-yes widely).
- *Hh04-g "Assault on properties", (Stealing, pillage, plundering )*, with two categories (1-no/ little, 2-yes widely).
- *Hh04-h "Physical violence", (Beating, rape, etc)*, with two categories (1-no/ little, 2-yes widely).

Notes: In this set of variables we considered the answer don't know as missing system, as this answer does not give an opinion

Another set of questions talked about the child himself:

- *Hh01-a "How do you evaluate your physical health status"* with two categories (1-good, 2-moderate/poor).
- *Hh01-b "How do you evaluate your mental health status"* with two categories (1-good, 2-moderate/poor).
- *Hr04 "Sex"* with two categories (1-male, 2-female).
- *Hh02 "Do you want your current weight to "*, we merged the categories to three only (1-remain as it is, 2-to decrease, 3-to increase).
- *Hr08 "Enrolled in education status"*, with two categories (1-currently enrolled, 2-not enrolled now).
- *S01 "Free time you have"* with three categories (1-little, 2-enough, 3- too much).

Another set of questions talked about the family circumstances:

- *Loctype "locality type"* with two categories (1-urban, 2- camps).
- *Ir04 "Total number of household members"* (numeric variable).
- *Hr07 "Refugee status"* with two categories (1-refugee, 2-not refugee).
- *Ho5 "Current status of parents"* the variable was with 7 categories (1-living together, 2-divorced, 3-father is dead, 4-mother is dead, 5-both are dead, 6-one of them works abroad, 7-others), some of these without frequencies, so we merged the categories to three only as: 1-living together, 2-one of them dead, 3-divorced & others.

We used the same variable's name and codes used by PCBS survey. Full detailed of the explanatory variables was summarized in table 3.2. This table prepared by using frequency command of (Statistical Package for Social Sciences) SPSS.



**Table 3.2: List of the explanatory variables and their frequencies**

| Explanatory variables |                        | Frequency     | %     | Valid %   | Cumulative % |
|-----------------------|------------------------|---------------|-------|-----------|--------------|
| Ho5                   | 1-living together      | 62,905        | 94.0  | 94.0      | 94.0         |
|                       | 2-one of them dead     | 2,615         | 3.9   | 3.9       | 97.9         |
|                       | 3-divorced 7 other     | 1,415         | 2.1   | 2.1       | 100.0        |
| S01                   | 1-little               | 14,134        | 21.1  | 21.1      | 21.1         |
|                       | 2-enough               | 32,903        | 49.2  | 49.2      | 70.3         |
|                       | 3-too much             | 19,898        | 29.7  | 29.7      | 100.0        |
| Hh01-a                | 1-good                 | 61,143        | 91.3  | 91.3      | 91.3         |
|                       | 2-moderate/poor        | 5,792         | 8.7   | 8.7       | 100.0        |
| Hh01-b                | 1-good                 | 53,884        | 80.4  | 80.4      | 80.4         |
|                       | 2-moderate/poor        | 13,091        | 19.6  | 19.6      | 100.0        |
| Hh02                  | 1-remain as it is      | 29,382        | 43.9  | 48.9      | 48.9         |
|                       | 2-to decrease          | 12,292        | 18.4  | 20.4      | 69.3         |
|                       | 3-to increase          | 18,446        | 27.6  | 30.7      | 100.0        |
|                       | 4-don't know (missing) | 6,815         | 10.2  |           |              |
| Hh04-a                | 1-n0/little            | 44,472        | 66.4  | 91.3      | 91.3         |
|                       | 2-yes widely           | 4,250         | 6.3   | 8.7       | 100.0        |
|                       | 3-don't know (missing) | 18,212        | 27.2  |           |              |
| Hh04-b                | 1-n0/little            | 10,058        | 15.0  | 15.2      | 15.2         |
|                       | 2-yes widely           | 55,982        | 83.6  | 84.8      | 100.0        |
|                       | 3-don't know (missing) | 895           | 1.3   |           |              |
| Hh04-c                | 1-n0/little            | 43,809        | 65.5  | 68.1      | 68.1         |
|                       | 2-yes widely           | 20,560        | 30.7  | 31.9      | 100.0        |
|                       | 3-don't know (missing) | 2,565         | 3.8   |           |              |
| Hh04-d                | 1-n0/little            | 41,617        | 62.2  | 91.4      | 91.4         |
|                       | 2-yes widely           | 3,920         | 5.8   | 8.6       | 100.0        |
|                       | 3-don't know (missing) | 21,398        | 32.0  |           |              |
| Hh04-e                | 1-n0/little            | 34,610        | 51.7  | 52.6      | 52.6         |
|                       | 2-yes widely           | 31,169        | 46.6  | 47.4      | 100.0        |
|                       | 3-don't know (missing) | 1,156         | 1.7   |           |              |
| Hh04-f                | 1-n0/little            | 50,580        | 75.6  | 79.6      | 79.6         |
|                       | 2-yes widely           | 12,952        | 19.4  | 20.4      | 100.0        |
|                       | 3-don't know (missing) | 3,403         | 5.1   |           |              |
| Hh04-g                | 1-n0/little            | 52,709        | 78.7  | 85.9      | 85.9         |
|                       | 2-yes widely           | 8,669         | 13.0  | 14.1      | 100.0        |
|                       | 3-don't know (missing) | 5,557         | 8.3   |           |              |
| Hh04-h                | 1-n0/little            | 48,037        | 71.8  | 77.7      | 77.7         |
|                       | 2-yes widely           | 13,813        | 20.6  | 22.3      | 100.0        |
|                       | 3-don't know (missing) | 5,084         | 7.6   |           |              |
| Hr04                  | 1-male                 | 32,387        | 48.4  | 48.4      | 48.4         |
|                       | 2-female               | 34,548        | 51.6  | 51.6      | 100.0        |
| Hr07                  | 1-refugee              | 39,840        | 59.5  | 59.5      | 59.5         |
|                       | 2-not refugee          | 27,095        | 40.5  | 40.5      | 100.0        |
| Hr08                  | Currently enrolled     | 66,374        | 99.2  | 99.2      | 99.2         |
|                       | Not enrolled now       | 560           | 0.8   | 0.8       | 100.0        |
| Loctype               | 1-urban                | 53,791        | 80.4  | 80.4      | 80.4         |
|                       | 2-camps                | 13,144        | 19.6  | 19.6      | 100.0        |
| Ir04                  | Mean= 9.28             | Std deviation | 2.529 | Median= 9 |              |

#### **IV. Building of MLR model**

We chose 18 explanatory variables which we believed it had an influence on physical violence issues against children see table 4.2. We tried to explore the effects of these variables by building MLR model and then examined of the results. To achieve this goal, we used SPSS software version 13, and used NOMERG command to calculate the MLR model with response variable and all explanatory variables to make the primary model.

#### **Checking of response variable**

From the table 4.1 of case processing summary we can check some points:

**Table 4.1: Case processing summary by using 18 explanatory variables**

| Response variable categories | N        | Marginal percentage |
|------------------------------|----------|---------------------|
| 0 Had not been               | 25317.33 | 75.7                |
| 1 Father/Mother              | 2858.29  | 8.6                 |
| 2 Sibling                    | 3235.41  | 9.7                 |
| 3 Peer & other               | 2011.97  | 6.0                 |
| Valid                        | 33423.00 | 100.0               |
| Missing                      | 33511.57 |                     |
| Total                        | 66934.56 |                     |
| Subpopulation                | 126      |                     |

*The dependent variable has only one value observed in 125(99.2%) subpopulation*

Table 4.1 is a portion of large table contains all variables, response variable and explanatory variables. We are focusing on the response variable, as we see in Table 4.3; the number of the valid observations used in our model is 33,423 distributed among the four categories. The marginal percentage column lists the proportion of valid observations found in each of the response variable' groups, 75.7% of the valid case (had not been) subjected to physical violence, 8.6% had been subjected by (father/mother), 9.7% by (sibling), and 6.0% by (peer and other).

#### **Subpopulation**

Subpopulation indicates the number of subpopulations contained in the data. A subpopulation of the data consists of one combination of the explanatory variables specified for the model. The SPSS footnote for table 4.1 provides how many of these combinations of the explanatory variables consist of records that all have the same value in the response variable. In our model there are 126 combinations that appear in the data and 125 of these combinations are composed of records with the same response variable categories

## Missing

Missing indicates the number of cases in the dataset where data are missing of the response variable or any of explanatory variables. In primary model we found the missing almost 50%. Brannon et al (2007) suggests that we can calculate scales with missing items if at least two thirds of the items were completed and others were dropped. Anyhow, this model still under checking, but we refer to some explanatory variables like alcohol consumption, smoking, etc, a category of "I don't know" is not a valid decision in this situation. It was considered as missing system as this procedure will not affect the final result, Moorman and Carr (2008).

## Computing by chance accuracy

The proportional by chance accuracy rate can be computed by calculating the proportion of cases for each group based on the number of cases in each group of the response variable. The squaring and summing the proportion of cases in each group are  $(0.757^2 + 0.086^2 + 0.097^2 + 0.06^2) = 0.593454 = 59.34\%$  (rounding error taken into account). The benchmark that used to characterize MLR model as useful is a 25% improvement over the rate of accuracy achievable by chance alone, so the proportional by chance accuracy criteria are:  $1.25 * 0.593454 = 74.18\%$ . This proportion will be compared with the overall percentage of the final model.

## Checking of explanatory variables

The ML method used to calculate MLR by using an iterative fitting process that attempts to cycle through repetitions to find an answer. Sometimes the method will break down and not be able to converge or find an answer, or sometimes will produce widely improbable results. Reporting that one- unit change in an explanatory variable increases the odds of the model by unreasonable results. These implausible results can be produced by multicollinearity, or categories of explanatory variables having no cases or zero cells. If we faced a situation comparable to that, this means we have a numerical problem and should not interpret the results. The practical solution of this problem is checking the standard errors for the parameter's explanatory variables that are larger than 2, Schwab (2007). This information is available in the parameter estimates table of the MLR results. This table 4.4 is consisting of three parts as the response variable has four categories. After checking this table we found there were five of the explanatory variables had standard error more than 2, we summarized the results of these variables in the Table 4.2.

**Table 4.2: The parameter estimates with more than 2 units of standard error**

| Response physical violence by | B       | Std .Error | Wald  | df | sig  | Exp(B)      |
|-------------------------------|---------|------------|-------|----|------|-------------|
| 1-father/mother               |         |            |       |    |      |             |
| H05=1                         | 14.510  | 160.461    | 0.008 | 1  | .928 | 2001804.550 |
| H05=2                         | -.652   | 190.610    | 0.0   | 1  | .997 | .521        |
| Hh04-b=1                      | -14.034 | 46.734     | 0.090 | 1  | .764 | 8.04E-007   |
| 2-sibling                     |         |            |       |    |      |             |
| H05=1                         | 14.050  | 152.557    | 0.008 | 1  | .927 | 1264463.302 |
| H05=2                         | 13.386  | 152.557    | 0.008 | 1  | .930 | 650756.508  |
| 3-peer & other                |         |            |       |    |      |             |
| H05=1                         | 12.302  | 103.476    | 0.014 | 1  | .905 | 220089.142  |
| H05=2                         | -5.735  | 135.662    | 0.002 | 1  | .966 | .003        |
| Hh01-a=1                      | 35.119  | 46.309     | 0.575 | 1  | .488 | 1786163.... |
| Hh02=2                        | -32.769 | 34.640     | 0.895 | 1  | .344 | 5.87E-015   |
| Hr08=1                        | 15.417  | 319.464    | 0.002 | 1  | .962 | 4960497.108 |

From Table 4.2, there are five explanatory variables causing a numerical problem, these variables: *H05 "current status of parents"*, *hh04-b "smoking among youth"*, *hh01-a "evaluation of physical health status"*, *hh02 "current weight"*, and *hr08 "enrolling in education"*. We note these parameter estimates gave unreasonable results by one unit change in the explanatory variable.

### Selection of the model

In the second phase we re-calculated the model after excluding the five variables, which had causing a numerical problem, and scanned the result and re-calculate of the model again after excluding the variables that had parameter estimates not statistically significant. This procedure had been replicated four times, and then we stopped as we found all parameter estimates are significant. The results were summarized for these four models in Table 4.3.

**Table 4.3: The specifications of four models**

| Description                                | Model (1) | Model(2) | Model(3) | Model(4) |
|--|-----------|----------|----------|----------|
| Number of explanatory variables            | 13        | 10       | 8        | 6        |
| Valid cases                                | 37202     | 37715    | 37715    | 38184    |
| Missing cases                              | 29732     | 29219    | 29219    | 28749    |
| Subpopulation                              | 133       | 117      | 82       | 31       |
| Chi-square value (likelihood ratio test)   | 13423.84  | 12425.41 | 9817.82  | 9073.22  |
| Df   | 42        | 33       | 24       | 18       |
| R- square Cox and Senell                   | 0.303     | 0.281    | 0.229    | .211     |
| R-square Nagelkerke                        | 0.377     | 0.351    | 0.287    | .265     |
| R-square Mc Fadden                         | 0.222     | 0.205    | 0.162    | .149     |
| Number of correct predicted "had not been  | 27473     | 28158    | 28180    | 28798    |
| Number of correct predicted "father/mother | 571       | 571      | 571      | 571      |
| Number of correct predicted "sibling"      | 342       | 1198     | 513      | 684      |
| Number of correct predicted "peer & other" | 520       | 520      | 0        | 0        |
| Classification Overall percentage          | 77.7      | 80.7     | 77.6     | 78.7     |

Information of the four models showed that the Model (2) is the best to be appropriate to the data comparing with the other models. It has the highest classification overall percentage, includes 10 independent variables. Also worked to increase the valid cases and reduce the missing cases. R-square factor which is usually influenced by a number of variables, had given the values comparable to the other models. For these reasons we selected model (2) with the following explanatory variables (10 variables): *Hh01-b* "Evaluation of mental health status", *Hh04-a* "Alcohol consumption", *Hh04-d* "Drug abuse", *Hh04-g* "Assault on properties", *Hh04-h* "Physical violence", *Hr04* "Sex", *Hh04-f* "Begging", *Ir04* "Total number of household members", *S01* "Free time you have", *Hh04-e* "Verbal violence".

### Checking of the selected model

There is no guarantee that the model fits the data well for any particular logistic regression model, Agresti (2007). Even though, there some ways used to detect lack of fit like likelihood ratio test.

### The sample size requirements

The minimum number of valid cases for each of explanatory variable according to a guideline provided by Hosmer and Lemeshow is 10, and preferring case to variable ratio 20 to 1 explanatory variable, Schwab (2007). In selected model the ratio is 377 cases to 1 explanatory variable. Table 4.4 presents the response variable and their categories frequencies used in selected model.

**Table 4.4: Case processing summary**

| Response variable Categories | N        | Marginal % |
|------------------------------|----------|------------|
| 0- Had not been              | 28767.61 | 76.3       |
| 1-Father/Mother              | 3091.03  | 8.2        |
| 2-Sibling                    | 3577.71  | 9.5        |
| 3-Peer & Other               | 2279.11  | 6.0        |
| Valid                        | 37715.46 |            |
| Missing                      | 29219.10 |            |
| Total                        | 66934.56 |            |
| Subpopulation                | 117      |            |

### Pseudo R-square

There are three pseudo R-square values can be calculated by SPSS for logistic regression table 4.5. Pseudo R –square does not have an equivalent to  $R^2$  in OLS regression (the coefficient of determination).  $R^2$  summarizes the proportion of variance in the response variable associated with explanatory variables, but pseudo R-square does not means what  $R^2$  means in OLS regression but we can use it as indicator for different areas of application. The model with the largest pseudo R-square statistic is best according to the measures; however,

classification coefficients as overall affect size measures are preferred over pseudo R-square measures as they have some severe limitations for this purpose, Garson (2009).

**Table 4.5: Pseudo R-Square**

|               |       |
|---------------|-------|
| Cox and Snell | 0.281 |
| Nagelkerke    | 0.351 |
| McFadden      | 0.205 |

### **Comparing accuracy rates**

The classification overall percentage computed by SPSS was 80.7% while was greater than the proportional by chance accuracy criterion of 74.1% or  $(1.25 \times 59.34) = 74.1$ .

The criterion for classification accuracy is satisfied. See Table 4.6.

**Table 4.6: The classification table of the selected model**

| Observed           | Predicted    |               |         |              |                 |
|--------------------|--------------|---------------|---------|--------------|-----------------|
|                    | Had not been | Father/Mother | Sibling | Peer & other | Percent correct |
| Had not been       | 28158.16     | 171.15        | 438.30  | 0            | 97.9%           |
| Father/Mother      | 2220.95      | 571.87        | 0       | 298.21       | 18.5%           |
| Sibling            | 2246.08      | 133.57        | 1198.05 | 0            | 33.5%           |
| Peer & other       | 1758.86      | 0             | 0       | 520.25       | 22.8%           |
| Overall percentage | 91.2         | 2.3           | 4.3     | 2.2          | 80.7%           |

### **Absence of Multicollinearity**

Multicollinearity can be occurred in logistic regression, as the correlation increases among the independent variables, the standard errors of the logit parameters will become inflated. Multicollinearity does not change the estimates of the parameters, only their reliability, Garson (2009). We check first the standard error; variables with more than 2 units were ignored. We checked the asymptotic correlation matrix which is a matrix of parameter estimate correlation. In this matrix we found that the majority of correlation coefficients were less than 0.10, another 4 were between (0.20 and 0.27), only one coefficient was 0.54, this means we do not have serious problem with multicollinearity among the explanatory variables that used in the model. Correlation between total number of household members (ir04) and evaluation of mental health status (hh01-b) is 0.221. Correlation between alcohol consumption among youth in locality (hh04-a) and drug abuse among youth in locality (hh04-d) is 0.208, also between drug abuse and assault on properties (hh04-g) is 0.54, between assault on properties and physical violence (hh04-h) is 0.276, and between the physical violence and begging (hh04-f) is 0.265.

### Goodness- of-fit measures

The likelihood ratio test is based on deviance  $[-2 \text{ Log Likelihood (LL)}]$ , the significance of the difference between the  $(-2LL)$  for our selected model minus likelihood ratio for a reduced model (intercept only) as in Table 4.9. A common use of the likelihood ratio test is to test this difference (it called *chi-square model*) dropping an interaction effect. If the chi-square model is significant, the interaction effect is contributing significantly to the full model and should be retained. The presence of a relationship between the response variable and combination of explanatory variables is based on the statistical significance of the final model chi-square. In our model, the p-value of the model chi-square (12425.41) was 0.000, less than the level of significance 0.05. We reject the null hypothesis which states that there was no difference between the model without explanatory variables and the model with explanatory variables. The existence of a relationship between the explanatory variables and the response variable was supported.

**Table 4.7: Model fitting information**

| Model          | Model Fitting Criteria |           |                  | Likelihood Ratio Tests |    |       |
|----------------|------------------------|-----------|------------------|------------------------|----|-------|
|                | AIC                    | BIC       | -2log Likelihood | Chi-Square             | df | Sig   |
| Intercept only | 53577.558              | 53603.171 | 53571.558        |                        |    |       |
| Final          | 41218.146              | 41525.508 | 41146.146        | 12425.411              | 33 | 0.000 |

AIC (Akaike Information Criterion) and BIC (Bayesian Information Criterion) judge a model by how close its fitted values tend to be to the true expected values, as summarized by a certain expected distance between the two, the optimal model is the one that tends to have its fitted values closest to the true outcome probabilities. In our model AIC and BIC and -2log likelihood are very close.

**Table (4.8): Likelihood ratio tests of the selected model**

| Effect    | Model Fitting Criteria |                      |                                    | Likelihood Ratio Tests |    |      |
|-----------|------------------------|----------------------|------------------------------------|------------------------|----|------|
|           | AIC of Reduced Model   | BIC of Reduced Model | -2 Log Likelihood of Reduced Model | Chi-Square             | df | Sig. |
| Intercept | 41218.146              | 41525.508            | 41146.146 <sup>a</sup>             | .000                   | 0  | .    |
| ir04      | 41434.453              | 41716.202            | 41368.453                          | 222.307                | 3  | .000 |
| hh01_b    | 42633.441              | 42915.189            | 42567.441                          | 1421.294               | 3  | .000 |
| hh04_a    | 42357.903              | 42639.651            | 42291.903                          | 1145.756               | 3  | .000 |
| hh04_d    | 42920.768              | 43202.517            | 42854.768                          | 1708.622               | 3  | .000 |
| hh04_g    | 42830.898              | 43112.646            | 42764.898                          | 1618.752               | 3  | .000 |
| hh04_h    | 43045.503              | 43327.251            | 42979.503                          | 1833.357               | 3  | .000 |
| hr04      | 44271.829              | 44553.577            | 44205.829                          | 3059.683               | 3  | .000 |
| hh04_f    | 41683.817              | 41965.565            | 41617.817                          | 471.670                | 3  | .000 |
| s01       | 42895.663              | 43151.798            | 42835.663                          | 1689.517               | 6  | .000 |
| hh04_e    | 42475.932              | 42757.680            | 42409.932                          | 1263.785               | 3  | .000 |

The chi-square statistic is the difference in -2 log-likelihoods between the final model and a reduced model.

The reduced model is formed by omitting an effect from the final model. The null hypothesis is that all parameters of that effect are 0.

- a. This reduced model is equivalent to the final model because omitting the effect does not increase the degrees of freedom.

In Table 4.8, we checked the same point with all explanatory variables used to build model separately. The result was referred that the existence of a relationship between each of the explanatory variables and the response variable was supported.

### **The criteria of Odds ratio explanation**

The major difference from the binomial situation is explanatory variable in selected model has three (b) parameters with three odds ratios (exp (b)), one for each level of the response physical violence except its reference category (0=had not been), (Table 4.9), which is assumed but does not show in the tables. The "exp(b)" column in SPSS's label for odds ratio of the explanatory variables with the response variable, it is predicted change in odds for a unit increase in the corresponding explanatory variable. Odds ratios less than 1 correspond to decreases and odds ratio more than 1.0 correspond to increases. Odds ratios close to 1.0 indicates that unit changes in that explanatory variable does not affect the response variable.

### **Estimating response probabilities**

The MLR model has an alternative expression in terms of the responses probabilities, that is  $\pi_j = \frac{e^{\alpha_j + \beta_j x}}{\sum_h e^{\alpha_h + \beta_h x}}$ ,  $j=1, \dots, J$ . In our model, we will denote the

probability of the child had not been faced physical violence (baseline category) by  $\pi_0$  and the estimate by  $\hat{\pi}_0$ . The physical violence by father/ mother by  $\pi_1$  and the estimate by  $\hat{\pi}_1$ . The physical violence by sibling by  $\pi_2$  and the estimate by  $\hat{\pi}_2$ , and the physical violence by peer and other by  $\pi_3$ , and the estimate by  $\hat{\pi}_3$ ,

the response probability satisfying  $\sum_{j=0}^3 \pi_j = 1$ , our baseline category is (had not been=0), from table (4.9) of parameter estimates we can calculate these probabilities by two steps:

First, we can calculate  $\log\left(\frac{\hat{\pi}_1}{\hat{\pi}_0}\right)$ ,  $\log\left(\frac{\hat{\pi}_2}{\hat{\pi}_0}\right)$ , and  $\log\left(\frac{\hat{\pi}_3}{\hat{\pi}_0}\right)$ , as the response variable has four categories (J=4), which means that there are 3 equations as following:

let  $y_1 = \log\left(\frac{\hat{\pi}_1}{\hat{\pi}_0}\right)$ , and  $y_2 = \log\left(\frac{\hat{\pi}_2}{\hat{\pi}_0}\right)$ , and  $y_3 = \log\left(\frac{\hat{\pi}_3}{\hat{\pi}_0}\right)$ , so



$$y_1 = 1.381 - 0.116(ir04) - 0.162(hh01-b=1) - 0.244(hh04-a=1) - 1.895(hh04-d=1) \\ + 0.465(hh04-g=1) - 0.445(hh04-h=1) - 1.163(hr04=1) - 0.599(hh04-f=1) \\ + 0.677(s01=1) + 1.291(s01=2) - 1.501(hh04-e=1) \quad (4.1)$$

$$y_2 = - 1.954 - 0.044(ir04) + 1.069(hh01-b=1) + 0.940(hh04-a=1) + 1.908(hh04-d=1) \\ - 2.326(hh04-g=1) - 1.620(hh04-h=1) - 1.964(hr04=1) + 0.843(hh04-f=1) \\ - 0.539(s01=2) \quad (4.2)$$

$$y_3 = 0.082(ir04) - 1.753(hh01-b=1) - 2.417(hh04-a=1) - 1.192(hh04-d=1) \\ - 0.862(hh04-g=1) + 1.983(hh04-h=1) + 1.195(hr04=1) - 0.536(hh04-f=1) \\ - 0.886(s01=1) - 1.853(s01=2) + 0.839(hh04-e=1) \quad (4.3)$$

We cannot make corresponding statement about variables "s01=1" and "hh04-e=1" in equation (2), and intercept in equation (3) as that odds ratios are non-significant. As Agresti (2007) says *"Statistical significance should not be the sole criterion for whether to include a term in a model. It is sensible to include a variable that is important for the purposes of the study and report its estimated effect even if it is not statistically significant. Keeping it in the model may help reduce bias in estimating effects of the other predictors and may make it possible to compare results with other studies where the effect is significant"*

Second we calculate  $\hat{\pi}_1, \hat{\pi}_2, \hat{\pi}_3, \hat{\pi}_0$ , as following, where exp or e = 2.71828 is the base of the system of natural logarithms:

$$\hat{\pi}_1 = \frac{\exp(y_1)}{1 + \exp(y_1) + \exp(y_2) + \exp(y_3)} \quad (4.4)$$

$$\hat{\pi}_2 = \frac{\exp(y_2)}{1 + \exp(y_1) + \exp(y_2) + \exp(y_3)} \quad (4.5)$$

$$\hat{\pi}_3 = \frac{\exp(y_3)}{1 + \exp(y_1) + \exp(y_2) + \exp(y_3)} \quad (4.6)$$

$$\hat{\pi}_0 = \frac{1}{1 + \exp(y_1) + \exp(y_2) + \exp(y_3)} \quad (4.7)$$

Where the (1) term in each denominator and in the numerator of  $\hat{\pi}_0$  represents  $\exp(\hat{\alpha}_0 + \hat{\beta}_0 x)$ , for  $\hat{\alpha}_0 = \hat{\beta}_0 = 0$ , Agresti (2007).

# An Application on Multinomial Logistic Regression Model

Table (4.9): The parameter estimates of the selected model

| response Physical violence by <sup>a</sup> |            | B              | Std. Error | Wald     | df | Sig. | Exp(B) |
|--|------------|----------------|------------|----------|----|------|--------|
| 1 Father/Mother                            | Intercept  | 1.381          | .133       | 108.121  | 1  | .000 |        |
|  | ir04       | -.116          | .010       | 123.252  | 1  | .000 | .891   |
|  | [hh01_b=1] | -.162          | .049       | 10.789   | 1  | .001 | .850   |
|  | [hh01_b=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_a=1] | -.244          | .064       | 14.439   | 1  | .000 | .784   |
|  | [hh04_a=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_d=1] | -1.895         | .068       | 780.216  | 1  | .000 | .150   |
|  | [hh04_d=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_g=1] | .465           | .077       | 36.714   | 1  | .000 | 1.592  |
|  | [hh04_g=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_h=1] | -.445          | .054       | 67.278   | 1  | .000 | .641   |
|  | [hh04_h=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hr04=1]   | -1.163         | .046       | 645.567  | 1  | .000 | .312   |
|  | [hr04=2]   | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_f=1] | -.599          | .053       | 129.126  | 1  | .000 | .549   |
|  | [hh04_f=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [s01=1]    | .677           | .066       | 106.005  | 1  | .000 | 1.967  |
|  | [s01=2]    | 1.291          | .060       | 455.276  | 1  | .000 | 3.635  |
|  | [s01=3]    | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_e=1] | -1.501         | .050       | 888.151  | 1  | .000 | .223   |
|  | [hh04_e=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
| 2 Sibling                                  | Intercept  | -1.954         | .174       | 126.245  | 1  | .000 |        |
|  | ir04       | -.044          | .010       | 17.415   | 1  | .000 | .957   |
|  | [hh01_b=1] | 1.069          | .060       | 314.681  | 1  | .000 | 2.913  |
|  | [hh01_b=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_a=1] | .940           | .090       | 110.122  | 1  | .000 | 2.561  |
|  | [hh04_a=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_d=1] | 1.908          | .096       | 396.509  | 1  | .000 | 6.737  |
|  | [hh04_d=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_g=1] | -2.326         | .061       | 1468.420 | 1  | .000 | .098   |
|  | [hh04_g=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_h=1] | -1.620         | .046       | 1215.457 | 1  | .000 | .198   |
|  | [hh04_h=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hr04=1]   | -1.964         | .050       | 1556.703 | 1  | .000 | .140   |
|  | [hr04=2]   | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_f=1] | .843           | .057       | 218.119  | 1  | .000 | 2.323  |
|  | [hh04_f=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [s01=1]    | -.020          | .049       | .169     | 1  | .681 | .980   |
|  | [s01=2]    | -.539          | .052       | 108.999  | 1  | .000 | .583   |
|  | [s01=3]    | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_e=1] | -.069          | .045       | 2.300    | 1  | .129 | .933   |
|  | [hh04_e=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
| 3 Peer& other                              | Intercept  | .145           | .146       | .980     | 1  | .322 |        |
|  | ir04       | .082           | .011       | 59.833   | 1  | .000 | 1.085  |
|  | [hh01_b=1] | -1.753         | .058       | 899.007  | 1  | .000 | .173   |
|  | [hh01_b=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_a=1] | -2.417         | .080       | 923.704  | 1  | .000 | .089   |
|  | [hh04_a=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_d=1] | -1.192         | .078       | 232.146  | 1  | .000 | .304   |
|  | [hh04_d=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_g=1] | -.862          | .077       | 126.518  | 1  | .000 | .422   |
|  | [hh04_g=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_h=1] | 1.983          | .100       | 395.992  | 1  | .000 | 7.262  |
|  | [hh04_h=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hr04=1]   | 1.195          | .064       | 347.519  | 1  | .000 | 3.305  |
|  | [hr04=2]   | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_f=1] | -.536          | .064       | 70.418   | 1  | .000 | .585   |
|  | [hh04_f=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [s01=1]    | -.886          | .063       | 199.821  | 1  | .000 | .412   |
|  | [s01=2]    | -1.853         | .068       | 733.513  | 1  | .000 | .157   |
|  | [s01=3]    | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |
|  | [hh04_e=1] | .839           | .061       | 192.266  | 1  | .000 | 2.315  |
|  | [hh04_e=2] | 0 <sup>b</sup> | .          | .        | 0  | .    | .      |

a. The reference category is: 0 Had not been.

b. This parameter is set to zero because it is redundant.

### Predications by using MLR model

Each case consists of a combination of explanatory variables. Prediction is based on classifying this combination in one of the four groups of the response variable. The model estimates the probabilities of this combination of the four groups of the response variable and then according to the largest probability will classify the case (we have 117 subpopulation group in the model). For the application of the model we had selected randomly of two cases of the data and we used the model to predict from which groups will classify it, the model consists of three equations, to estimate the four response probabilities  $(\pi_0, \pi_1, \pi_2, \pi_3)$  by using equations (4.1), (4.2), (4.3).

**Table 4.10: Selected cases information**

| Explanatory variables  | Case 12049     | Case 13102     |
|--|----------------|----------------|
| <i>Total number of household members (ir 04)</i>             | 8              | 4              |
| <i>Evaluation of the mental health status (hh01-b)</i>       | Moderate/ poor | Moderate/ poor |
| <i>Alcohol consumption among youth in locality (hh04-a)</i>  | No/little      | No/little      |
| <i>Drug abuse among youth in locality (hh04-d)</i>           | No/little      | Yes widely     |
| <i>Assault on properties among youth in locality(hh04-g)</i> | Yes widely     | Yes widely     |
| <i>Physical violence among youth in locality (hh04-h)</i>    | Yes widely     | Yes widely     |
| <i>Sex (hr04)</i>  | female         | Female         |
| <i>Begging among youth in locality (hh04-f)</i>              | Yes widely     | Yes widely     |
| <i>Free time child has (s01)</i>                             | enough         | Little         |
| <i>Verbal violence among youth in locality (hh04-e)</i>      | Yes widely     | Yes widely     |

#### For case number 12049,

By using information from table (4.9)

$$y_1 = \log \left( \frac{\hat{\pi}_1}{\hat{\pi}_0} \right) = 1.381 - 0.116(8) - 0.162(0) - 0.244(1) - 1.895(1) + 0.465(0) - 0.445(0) - .163(0) - 0.599(0) + 0.677(0) + 1.291(1) - 1.501(0) = \mathbf{-0.395}$$

$$y_2 = \log \left( \frac{\hat{\pi}_2}{\hat{\pi}_0} \right) = -1.954 - 0.044(8) + 1.069(0) + 0.940(1) + 1.908(1) - 2.326(0) - 1.620(0) - 1.964(0) + 0.843(0) - 0.539(0) = \mathbf{0.542}$$

$$y_3 = \log \left( \frac{\hat{\pi}_3}{\hat{\pi}_0} \right) = +0.082(8) - 1.753(0) - 2.417(1) - 1.192(1) - 0.862(0) + 1.983(0) + 1.195(0) - 0.536(0) - 0.886(0) - 1.853(1) + 0.839(0) = \mathbf{-4.806}$$

By using equations (4.4, 4.5, 4.6, and 4.7), we can calculate the estimated probability to occur in each category as the following:

$$\hat{\pi}_1 = \frac{\exp(-0.395)}{1 + \exp(-0.395) + \exp(0.542) + \exp(-4.806)} = 0.1980$$

$$\hat{\pi}_2 = \frac{\exp(0.542)}{1 + \exp(-0.395) + \exp(0.542) + \exp(-4.806)} = 0.5055$$

$$\hat{\pi}_3 = \frac{\exp(-4.806)}{1 + \exp(-0.395) + \exp(0.542) + \exp(-4.806)} = 0.0024$$

$$\hat{\pi}_4 = \frac{1}{1 + \exp(-0.395) + \exp(0.542) + \exp(-4.806)} = 0.2940$$

These probabilities appeared that the case number 12049 has probability of 0.198 to occur in category that the child had facing physical violence by his father/mother, and probability of 0.5055 by sibling, and 0.0024 by peer and other and finally had not been facing physical violence with probability of 0.2940. So the conclusion here is that the child was facing physical violence by sibling has the largest probability comparing with other groups or categories.

**For case number 13102,**

$$y_1 = \log\left(\frac{\hat{\pi}_1}{\hat{\pi}_0}\right) = 1.381 - 0.116(4) - 0.162(0) - 0.244(1) - 1.895(0) + 0.465(0) - 0.445(0) - 1.163(0) - 0.599(0) + 0.677(1) + 1.291(0) - 1.501(0) = \mathbf{1.35}$$

$$y_2 = \log\left(\frac{\hat{\pi}_2}{\hat{\pi}_0}\right) = -1.954 - 0.044(4) + 1.069(0) + 0.940(1) + 1.908(0) - 2.326(0) - 1.620(0) - 1.964(0) + 0.843(0) - 0.539(1) = \mathbf{-1.729}$$

$$y_3 = \log\left(\frac{\hat{\pi}_3}{\hat{\pi}_0}\right) = +0.082(4) - 1.753(0) - 2.417(1) - 1.192(0) - 0.862(0) + 1.983(0) + 1.195(0) - 0.536(0) - 0.886(1) - 1.853(0) + 0.839(0) = \mathbf{-2.975}$$

By using equations (4.4, 4.5, 4.6, and 4.7), we can calculate the estimated probability to occur in each category as the following:

$$\hat{\pi}_1 = \frac{\exp(1.35)}{1 + \exp(1.35) + \exp(-1.729) + \exp(-2.975)} = \mathbf{0.7584}$$

$$\hat{\pi}_2 = \frac{\exp(-1.729)}{1 + \exp(1.35) + \exp(-1.729) + \exp(-2.975)} = \mathbf{0.0349}$$

$$\hat{\pi}_3 = \frac{\exp(-2.975)}{1 + \exp(1.35) + \exp(-1.729) + \exp(-2.975)} = \mathbf{0.0100}$$

$$\hat{\pi}_4 = \frac{1}{1 + \exp(1.35) + \exp(-1.729) + \exp(-2.975)} = \mathbf{0.1966}$$

These probabilities appeared that the case number 13102 has probability of 0.7584 to occur in category that the child had facing physical violence by his father/mother, and probability of 0.0349 by sibling, and 0.0100 by peer and other and finally had not been facing physical violence with probability of 0.1966. So the conclusion here is that the child was facing physical violence by father/mother has the largest probability comparing with other groups or categories.

## **V. Conclusion**

We have reviewed the results of the model and carried out some tests to make sure that the model is fit of the data according to statistical terms. Also we have reviewed the estimates of parameters and interpreted these estimates focusing on odds ratio scale. Likelihood ratio tests showed all explanatory variables were significance but the effects and contribution of each variable were not the same, so it were sorted according to their effects on the model. "Sex" variable was the most significant, followed by "Spread of physical violence among young people", "Spread of drug abuse", "Free time the children had", "Assault on properties among the youth", "Mental status", "Spread of verbal violence", "Spread of alcohol among youth", "Spread of begging", and finally "Total number of household members". The model ability of prediction had been checked by choosing two cases of the data randomly and applying the model to predict in any of the response variable's group can be classified of these cases. The model has been successful in one classification.

The crucial conclusion can be presented by several important points:

1. The usage of the MLR model gives us the opportunity to deal with a response categorical variable with more than two levels and variety of explanatory variables.
2. MLR indicates the effect of each of explanatory variables as well as its additive effect by used in the analysis simultaneously which we are aiming of the study of this model.
3. MLR enables building a statistical model showing those complex and interrelated relationships, particularly as we are dealing with a qualitative response variable has more than two categories. These equations could measure accurately the effect of each of explanatory variables and excluded those variables which did not have statistical significant.
4. MLR model, also has proved its ability to predict, and has reached the precision with which exhibited 80.7% in our model.
5. The model will help researcher who will try to study the subject of physical violence by gave him an idea about variables importance and effects, of course it can be made comparisons between the effects that are calculated from models if used the similar variables.
6. The logistic regression model is a suitable model to many types of data when the response variable with more than two categories. MLR has no any restrictions about the explanatory variables; this model is most common in the categorical data analysis. MLR can be used in many areas of social, educational, health, behavioral and even scientific experiments.

## References

1. Agresti, A. (2007). *An Introduction to Categorical Data Analysis*. John Wiley & Sons, Inc.
2. Agresti, A. (2002). *Categorical Data Analysis*. John Wiley & Sons, Inc.
3. Brannon, D., Barry, T., Kemper, P., Schreiner, A., and Vasey, J.(2007). Job Perception and Intent to Leave Among Direct Care Workers: Evidence from the Better Jobs Better Care Demonstrations. *The Gerontologist* 47:820-829, The Gerontological Society of America.
4. Chatterjee, S., and Hadi, A. (2006). *Regression Analysis by Example*. John Wiley & Sons.
5. Garson, D. (2009). *Logistic Regression with SPSS*. North Carolina State University, Public administration Program.
6. Moorman, S.M., and Carr, D. (2008). Spouses Effectiveness as End-of-Life Health Care Surrogates: Accuracy, Uncertainty, and Errors of Overtreatment or Undertreatment. *The Gerontologist* 48:811, The Gerontological Society of America.
7. PCBS (2003). *User guide, Youth Survey 2003*. Palestinian Central Bureau of Statistics.
8. Schafer J.L. (2006). Multinomial logistic regression models. *STAT* 544-Lecture 19.
9. Schwab J., (2007). Multinomial Logistic Regression Basic Relationships. [www.utexas.edu/courses/schwab/sw388r7/SolvingProblems/Analyzi](http://www.utexas.edu/courses/schwab/sw388r7/SolvingProblems/Analyzi).
10. Takagi, E., Silverstein, M., and Crimmins, E. (2007). Intergenerational Coresidence of Older Adults in Japan: Conditions for Cultural Plasticity. *The Journals of Gerontology Series B: Psychological Sciences and Social Sciences* 62:S330-S339, The Gerontological Society of America.