

Hybrid Logistic Regression and Random-Response Model in a Matched Pairs Study

Chien-Hua Wu *

Department of Applied Mathematics
Chung-Yuan Christian University

Shu-Mei Wan

Lunghwa University of Science and Technology

Mei-Chi Li

Department of Applied Mathematics
Chung-Yuan Christian University

October 23, 2008

ABSTRACT

The randomized response technique is a survey procedure that respondent answers sensitive questions randomly. This article develops and illustrates the application of a covariate extension of the RRT for a matched pair data and that allows for modeling the relation between the proportion with a sensitive characteristic and a covariate.

1. INTRODUCTION

In 1965, Warner developed an interviewing procedure designed to reduce errors caused by nonresponse and untruthful answers of a sensitive question. Greeber, Abul-Ela, Simmons, and Hovitz (1969) modified the Warner model allowing the interviewer to ask question requiring a quantitative response to an unrelated question. Abul-Ela, Greenber, and Horvitz (1967) extended the randomized response technique to estimate multinomial proportions. Greenber, Kuebler, Abernathy, and Horvitz (1971) studied the randomized response technique in quantitative data. Sheers, Dayton, and Mitchell (1988) developed and illustrated the application of a covariate extension of the randomized response technique.

The purpose of this article is to introduce the randomized response technique in a matched pair study and to present the mathematical development of baseline categorical logit model for a covariate extension of the randomized response technique in matched pairs study. An example of 319 presidential gun shot, one day before 2004 presidential election in Taiwan, will be illustrated the use of the proposed method to see whether or not the 319 presidential gun shot alters the election result.

2. RRT OF MATCHED PAIRS DATA

The randomized response technique is a data collection procedure that allows researchers to obtain

*Corresponding author: Chien-Hua Wu, Tel:886-3-265-3130; Fax:886-3-265-3199. E-mail address: cwu@cycu.edu.tw

sensitive information while persons being interviewed often refuse to answer or give correct answers that may embarrass them or be harmful to them in some way. For example:

First Question:
Card A: *Did you plan to vote for candidate number one before 319 presidential gun shot?*
 Second Question:
Did you change you mind after 319 gun shot?

First Question:
Card B: *Did you plan to vote for candidate number two before 319 presidential gun shot?*
 Second Question:
Did you change your mind after 319 gun shot?

The respondent is directed to answer both questions in either Card A or Card B. The interviewer never knows which card is chosen. Let θ be the probability of Card A is selected. To develop an estimator of matched pair probability, we can think of the procedure outlined above as consisting of two stages: (1) select a card. (2) answer two questions in the selected card. This process can be modeled by a tree diagram. In Figure 1, parenthesis indicates his/her favor candidate number.

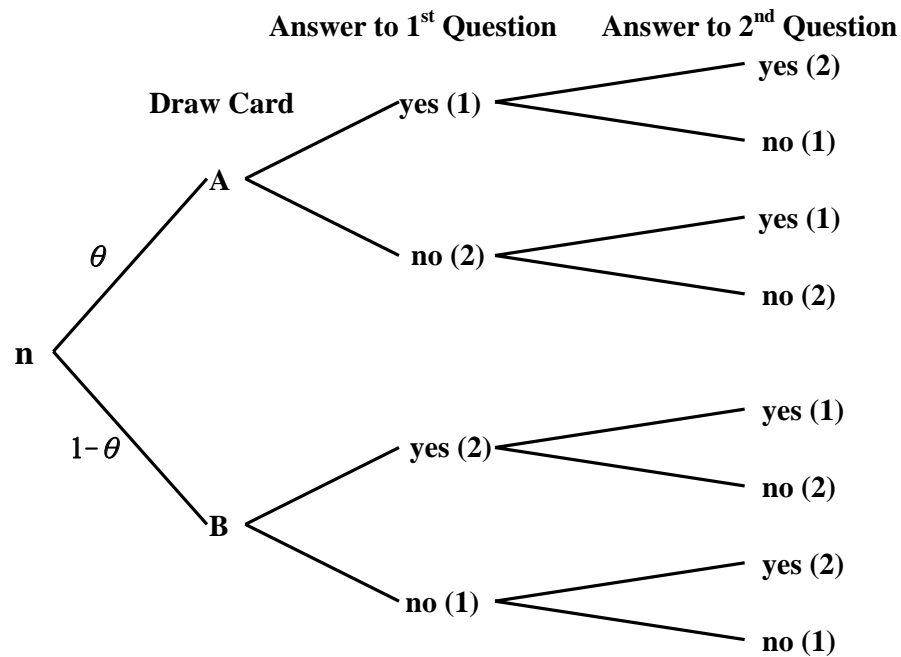


Figure 1: A Tree Diagram of Random-Response Model For A Matched Pair Data

A single random sample of n people is selected from the population. Each person in the sample is asked to randomly draw a card from the deck and to state "yes" if the question on the card agrees with the group to which he or she belongs, or "no" if the question on the card is different from the group to which he or she belongs. A two-way table having the same categories ("yes" or "no") for both questions summarizes such data. In table, the row marginal counts (n_{1+}, n_{2+}) are the numbers of totals for the first survey, and the column marginal counts (n_{+1}, n_{+2}) summarize results for the second survey. Let $\pi_{ij} = P(X = i, Y = j)$ denote the probability that (X, Y) falls in the cell in row i and column j , where $\sum_{i=1}^2 \sum_{j=1}^2 \pi_{ij} = 1$. The cell counts are denoted by $\{n_{ij}\}$, with $n = \sum_{i=1}^2 \sum_{j=1}^2 n_{ij}$. The corresponding sample joint proportion is $\hat{\pi}_{ij} = \frac{n_{ij}}{n}$, where $i, j = 1, 2$. The purpose of this study is to estimate the probability that a subject favors candidate number i before 319 presidential gun shot and votes for candidate number j after 319 presidential gun shot. Let p_{ij} denote such probability.

In the first survey question, there are two ways for the interviewer to obtain a "yes". Thus,

$$\begin{aligned} P(\text{yes}) &= \theta p_{1+} + (1 - \theta) p_{2+} \equiv \pi_{1+} \\ &\Rightarrow \hat{\pi}_{1+} = (2\theta - 1) \hat{p}_{1+} + (1 - \theta). \end{aligned}$$

Then the estimator of the probability of favor candidate number one before 319 gun shot is $\hat{p}_{1+} = \frac{\hat{\pi}_{1+} - (1 - \theta)}{2\theta - 1}$, where $\hat{\pi}_{1+} = \frac{n_{1+}}{n}$. Similarly, the probability of favor candidate number two before 319 gun shot can be estimated by $\hat{p}_{2+} = \frac{\theta - \hat{\pi}_{1+}}{2\theta - 1}$. Next, we consider to estimate the joint probabilities p_{ij} 's. Based on both questions, the four probabilities answering first question and second question can be written by

$$\begin{aligned} P(\text{yes}, \text{yes}) &= \theta p_{12} + (1 - \theta) p_{21} \equiv \pi_{11} \\ P(\text{yes}, \text{no}) &= \theta p_{11} + (1 - \theta) p_{22} \equiv \pi_{12} \\ P(\text{no}, \text{yes}) &= \theta p_{21} + (1 - \theta) p_{12} \equiv \pi_{21} \\ P(\text{no}, \text{no}) &= \theta p_{22} + (1 - \theta) p_{11} \equiv \pi_{22}. \end{aligned}$$

After some algebra, we can estimate the p_{ij} 's by

$$\begin{cases} \hat{p}_{22} = \frac{\theta \hat{\pi}_{+2} - \hat{\pi}_{12}}{2\theta - 1} \\ \hat{p}_{11} = \frac{\theta \hat{\pi}_{+2} - \hat{\pi}_{22}}{2\theta - 1} \\ \hat{p}_{21} = \frac{\theta \hat{\pi}_{+1} - \hat{\pi}_{11}}{2\theta - 1} \\ \hat{p}_{12} = \frac{\theta \hat{\pi}_{+1} - \hat{\pi}_{21}}{2\theta - 1} \end{cases}$$

, where $\hat{\pi}_{1+} = \frac{n_{1+}}{n}$ and $\hat{\pi}_{2+} = \frac{n_{2+}}{n}$.

3. HYPOTHESIS TESTING

The delta method for random vectors implies asymptotic normality of a function of cell counts in contingency table. The cell counts $(n_{11}, n_{12}, n_{21}, n_{22})$ have a multinomial distribution with cell probabilities $\boldsymbol{\pi} = (\pi_{11}, \pi_{12}, \pi_{21}, \pi_{22})'$. The multivariate Central Limit Theorem (Rao 1973) implies

$$\hat{\boldsymbol{\pi}} = (\hat{\pi}_{11}, \hat{\pi}_{12}, \hat{\pi}_{21}, \hat{\pi}_{22})' \xrightarrow{d} N\left(\boldsymbol{\pi}, \frac{\boldsymbol{\Sigma}}{n}\right)$$

, where $\Sigma = \text{diag}(\pi) - \pi\pi'$.

By the delta method, $\hat{p}_{12} - \hat{p}_{21}$, a function of $\hat{\pi}$, having nonzero differential at π are also asymptotically normally distributed.

$$(\hat{p}_{12} - \hat{p}_{21}) \xrightarrow{d} N \left((p_{12} - p_{21}), \frac{\phi' \Sigma \phi}{n} \right), \text{ where } \phi = (\phi_{11}, \phi_{12}, \phi_{21}, \phi_{22})' \text{ and}$$

$$\begin{aligned} \phi_{11} &= \left. \frac{\partial(\hat{p}_{12} - \hat{p}_{21})}{\partial \hat{\pi}_{11}} \right|_{\hat{\pi}=\pi} = \frac{1}{2\theta - 1} \\ \phi_{12} &= \left. \frac{\partial(\hat{p}_{12} - \hat{p}_{21})}{\partial \hat{\pi}_{12}} \right|_{\hat{\pi}=\pi} = 0 \\ \phi_{21} &= \left. \frac{\partial(\hat{p}_{12} - \hat{p}_{21})}{\partial \hat{\pi}_{21}} \right|_{\hat{\pi}=\pi} = \frac{-1}{2\theta - 1} \\ \phi_{22} &= \left. \frac{\partial(\hat{p}_{12} - \hat{p}_{21})}{\partial \hat{\pi}_{22}} \right|_{\hat{\pi}=\pi} = 0. \end{aligned}$$

Since $\phi' \Sigma \phi$ are continuous at π , $\hat{\phi}' \hat{\Sigma} \hat{\phi}$ is a consistent estimator of $\phi' \Sigma \phi$. Thus, confidence interval and test use the result that $\frac{\sqrt{n}[(\hat{p}_{12} - \hat{p}_{21}) - (p_{12} - p_{21})]}{\sqrt{\hat{\phi}' \hat{\Sigma} \hat{\phi}}}$ is asymptotically standard normal. For instance, the test statistic $\frac{\sqrt{n}(\hat{p}_{12} - \hat{p}_{21})}{\sqrt{\hat{\phi}' \hat{\Sigma} \hat{\phi}}}$ is a test of marginal homogeneity for matched binary responses has null hypothesis $H_0 G p_{12} = p_{21}$.

4. BASELINE CATEGORICAL LOGIT MODEL

In this section, we use baseline categorical logit model to describe effects of the explanatory variables, x , on matched pair data. Considering $p_{22}(x)$ as baseline, the logit pairs are given by

$$\begin{aligned} \log \left(\frac{p_{11}(x)}{p_{22}(x)} \right) &= \alpha_{11} + \beta_{11}x \\ \log \left(\frac{p_{12}(x)}{p_{22}(x)} \right) &= \alpha_{12} + \beta_{12}x \\ \log \left(\frac{p_{21}(x)}{p_{22}(x)} \right) &= \alpha_{21} + \beta_{21}x. \end{aligned}$$

An alternative formula for logistic regression refers directly to the success probability as follows.

$$\begin{aligned} p_{11}(x) &= \frac{\exp(\alpha_{11} + \beta_{11}x)}{1 + \exp(\alpha_{11} + \beta_{11}x) + \exp(\alpha_{12} + \beta_{12}x) + \exp(\alpha_{21} + \beta_{21}x)} \\ p_{12}(x) &= \frac{\exp(\alpha_{12} + \beta_{12}x)}{1 + \exp(\alpha_{11} + \beta_{11}x) + \exp(\alpha_{12} + \beta_{12}x) + \exp(\alpha_{21} + \beta_{21}x)} \\ p_{21}(x) &= \frac{\exp(\alpha_{21} + \beta_{21}x)}{1 + \exp(\alpha_{11} + \beta_{11}x) + \exp(\alpha_{12} + \beta_{12}x) + \exp(\alpha_{21} + \beta_{21}x)} \\ p_{22}(x) &= \frac{1}{1 + \exp(\alpha_{11} + \beta_{11}x) + \exp(\alpha_{12} + \beta_{12}x) + \exp(\alpha_{21} + \beta_{21}x)} \end{aligned}$$

Clearly, $\mathbf{n} = (n_{11}, n_{12}, n_{21}, n_{22})'$ follows a multinomial distribution with parameters $\boldsymbol{\pi} = (\pi_{11}, \pi_{12}, \pi_{21}, \pi_{22})'$. Let n_{lmi} denote the observed frequency of cell (l, m) at level i of x , where $l, m = 1, 2$. The likelihood function can be written by

$$L(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \prod_{i=1}^k \frac{n_i!}{n_{11i}!n_{12i}!n_{21i}!n_{22i}!} (\pi_{11}(x_i))^{n_{11i}} (\pi_{12}(x_i))^{n_{12i}} (\pi_{21}(x_i))^{n_{21i}} (\pi_{22}(x_i))^{n_{22i}}$$

, where $\boldsymbol{\alpha} = (\alpha_{11}, \alpha_{12}, \alpha_{21})'$, $\boldsymbol{\beta} = (\beta_{11}, \beta_{12}, \beta_{21})'$ and

$$\begin{aligned} \pi_{11}(x) &= \theta p_{12}(x) + (1 - \theta)p_{21}(x) \\ &= \frac{\theta \exp(\alpha_{12} + \beta_{12}x) + (1 - \theta)\exp(\alpha_{21} + \beta_{21}x)}{1 + \exp(\alpha_{11} + \beta_{11}x) + \exp(\alpha_{12} + \beta_{12}x) + \exp(\alpha_{21} + \beta_{21}x)} \\ \pi_{12}(x) &= \frac{\theta \exp(\alpha_{11} + \beta_{11}x) + (1 - \theta)}{1 + \exp(\alpha_{11} + \beta_{11}x) + \exp(\alpha_{12} + \beta_{12}x) + \exp(\alpha_{21} + \beta_{21}x)} \\ \pi_{21}(x) &= \frac{\theta \exp(\alpha_{21} + \beta_{21}x) + (1 - \theta)\exp(\alpha_{12} + \beta_{12}x)}{1 + \exp(\alpha_{11} + \beta_{11}x) + \exp(\alpha_{12} + \beta_{12}x) + \exp(\alpha_{21} + \beta_{21}x)} \\ \pi_{22}(x) &= \frac{\theta + (1 - \theta)\exp(\alpha_{11} + \beta_{11}x)}{1 + \exp(\alpha_{11} + \beta_{11}x) + \exp(\alpha_{12} + \beta_{12}x) + \exp(\alpha_{21} + \beta_{21}x)}. \end{aligned}$$

The equations that determine the maximal likelihood (ML) parameter estimates are obviously non-linear, and the estimates do not have a closed-form expression. To calculate the ML estimates of model parameters, a popular algorithm for doing this, called Newton-Raphson algorithm, is applied to approximate the log-likelihood function.

5. AN EXAMPLE OF 319 PRESIDENTIAL GUN SHOT

The Taiwan presidential election of 2004, March 20, was won by the incumbent President who defeated his main rival by a margin of 0.22% of valid votes. On March 19, one day before the presidential election, President and Vice-President were both shot while campaigning in Tainan, Taiwan. Their injuries were not life-threatening, and both President and Vice-President were released from hospital on the same day without losing consciousness or having surgery. The election campaigns on both sides were suspended, but the next day's election was not postponed, as Taiwanese law only allows for suspension of election upon the death of a candidate. After all 13,749 polling places had reported, the rival of the incumbent President had refused to concede and challenged the result. The main goal of this article is to investigate if 319 presidential gun shot alters the election result.

In Taiwan, persons being interviewed often refuse to answer or give correct answers to a political question. Thus, the development of the randomized response technique as a survey technique to eliminate evasive answer bias is necessary. The questionnaire has been designed in Section 2. If the last digit of the respondent's social security number falls in 1, 3 or 9, the interviewed respondent answers the questions in Card B. Otherwise, the respondent answers the those questions in Card A. They are to answer only "yes" or "no" for both questions. Based on a random sample of 500 voting-age Taiwan citizens, the cell counts and marginal totals are summarized in Table 1. For instance, the row marginal counts (243, 257) are the ("yes", "no") totals for the number of subjects answering the first question regardless from Card A or Card B, and the column marginal counts (113, 387) summarize results for the second question. The corresponding joint proportions and marginal proportions are estimated by

$$\hat{\pi}_{11} = \frac{n_{11}}{n} = \frac{69}{500} = 0.138, \hat{\pi}_{12} = \frac{n_{12}}{n} = \frac{174}{500} = 0.348$$

I \ II	yes	no	
yes	69	174	243
no	44	213	257
	113	387	500

Table 1: Observed Cell Counts of Answering Two Questions

$$\begin{aligned} \hat{\pi}_{21} &= \frac{n_{21}}{n} = \frac{44}{500} = 0.088 & \hat{\pi}_{22} &= \frac{n_{22}}{n} = \frac{213}{500} = 0.426 \\ \hat{\pi}_{1+} &= \frac{n_{1+}}{n} = \frac{243}{500} = 0.486 & \hat{\pi}_{2+} &= \frac{n_{2+}}{n} = \frac{257}{500} = 0.514 \\ \hat{\pi}_{+1} &= \frac{n_{+1}}{n} = \frac{113}{500} = 0.226 & \hat{\pi}_{+2} &= \frac{n_{+2}}{n} = \frac{387}{500} = 0.774 \end{aligned}$$

The major purpose of this study is to estimate the probability that a subject favors candidate number i before 319 presidential gun shot and votes for candidate number j after 319 presidential gun shot. The cell proportions of matched pairs in the randomized response technique can be obtained by

$$\left\{ \begin{aligned} \hat{p}_{11} &= \frac{\theta \hat{\pi}_{+2} - \hat{\pi}_{22}}{2\theta - 1} = \frac{0.7 \times 0.774 - 0.426}{2 \times 0.7 - 1} = 0.2895 \\ \hat{p}_{12} &= \frac{\theta \hat{\pi}_{+1} - \hat{\pi}_{21}}{2\theta - 1} = \frac{0.7 \times 0.226 - 0.088}{2 \times 0.7 - 1} = 0.1755 \\ \hat{p}_{21} &= \frac{\theta \hat{\pi}_{+1} - \hat{\pi}_{11}}{2\theta - 1} = \frac{0.7 \times 0.226 - 0.138}{2 \times 0.7 - 1} = 0.0505 \\ \hat{p}_{22} &= \frac{\theta \hat{\pi}_{+2} - \hat{\pi}_{12}}{2\theta - 1} = \frac{0.7 \times 0.774 - 0.348}{2 \times 0.7 - 1} = 0.4845. \end{aligned} \right.$$

Consequently, the another question of interest for such data is to test the marginal homogeneity for matched binary responses in RRT. Since

$$\begin{aligned} \phi_{11} &= \left. \frac{\partial(\hat{p}_{12} - \hat{p}_{21})}{\partial \hat{\pi}_{11}} \right|_{\hat{\pi}=\pi} = \frac{1}{2\theta - 1} = \frac{1}{2 \times 0.7 - 1} = 2.5 \\ \phi_{12} &= \left. \frac{\partial(\hat{p}_{12} - \hat{p}_{21})}{\partial \hat{\pi}_{12}} \right|_{\hat{\pi}=\pi} = 0 \\ \phi_{21} &= \left. \frac{\partial(\hat{p}_{12} - \hat{p}_{21})}{\partial \hat{\pi}_{21}} \right|_{\hat{\pi}=\pi} = \frac{-1}{2\theta - 1} = \frac{-1}{2 \times 0.7 - 1} = -2.5 \\ \phi_{22} &= \left. \frac{\partial(\hat{p}_{12} - \hat{p}_{21})}{\partial \hat{\pi}_{22}} \right|_{\hat{\pi}=\pi} = 0 \end{aligned}$$

The asymptotic variance equals

$$\begin{aligned}\phi' \hat{\Sigma} \phi &= \sum_{i=1}^2 \sum_{j=1}^2 \hat{\pi}_{ij} \phi_{ij}^2 - \left(\sum_{i=1}^2 \sum_{j=1}^2 \hat{\pi}_{ij} \phi_{ij} \right)^2 \\ &= \hat{\pi}_{11} \phi_{11}^2 + \hat{\pi}_{12} \phi_{12}^2 + \hat{\pi}_{21} \phi_{21}^2 + \hat{\pi}_{22} \phi_{22}^2 - (\hat{\pi}_{11} \phi_{11} + \hat{\pi}_{12} \phi_{12} + \hat{\pi}_{21} \phi_{21} + \hat{\pi}_{22} \phi_{22})^2 \\ &= 0.138 \times (2.5)^2 + 0.088 \times (-2.5)^2 - (0.138 \times 2.5 + 0.088 \times (-2.5))^2 \\ &= 1.396875\end{aligned}$$

Thus the test statistic is

$$\begin{aligned}Z &= \frac{\sqrt{n}(\hat{p}_{12} - \hat{p}_{21})}{\sqrt{\phi' \hat{\Sigma} \phi}} \\ &= \frac{0.125\sqrt{500}}{\sqrt{1.396875}} \\ &= 2.3649\end{aligned}$$

Comparing with $Z_{\alpha/2} = Z_{0.025} = 1.96$, we reject the null hypothesis $H_0 Gp_{12} = p_{21}$. There is evidence that the 319 presidential gun shot alters the election result.

Finally, the relation between gender and the cell proportions of matched pairs in RRT can be assessed by baseline categorical logit model. Since \mathbf{n} follows a multinomial distribution with parameters $\boldsymbol{\pi}$, the maximum likelihood function can be written by

$$L(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \prod_{i=1}^2 \frac{n!}{n_{11i}! n_{12i}! n_{21i}! n_{22i}!} (\pi_{11}(x_i))^{n_{11i}} (\pi_{12}(x_i))^{n_{12i}} (\pi_{21}(x_i))^{n_{21i}} (\pi_{22}(x_i))^{n_{22i}}$$

, where x_1 stands for female and x_2 represents for male. Applying Newton-Raphson algorithm, its ML estimates equal $\alpha_{11} = 0.0016928$, $\alpha_{12} = 15.251517$, $\alpha_{21} = 15.257606$, $\beta_{11} = -0.00045$, $\beta_{12} = -16.27752$, and $\beta_{21} = -16.28102$. Thus, the sample joint proportions for female are

$$\begin{aligned}p_{11}(x_1) &= \frac{\exp(\alpha_{11})}{1 + \exp(\alpha_{11}) + \exp(\alpha_{12}) + \exp(\alpha_{21})} = 0.0000001187766792 \\ p_{12}(x_1) &= \frac{\exp(\alpha_{12})}{1 + \exp(\alpha_{11}) + \exp(\alpha_{12}) + \exp(\alpha_{21})} = 0.4984776364 \\ p_{21}(x_1) &= \frac{\exp(\alpha_{21})}{1 + \exp(\alpha_{11}) + \exp(\alpha_{12}) + \exp(\alpha_{21})} = 0.5015221262 \\ p_{22}(x_1) &= \frac{1}{1 + \exp(\alpha_{11}) + \exp(\alpha_{12}) + \exp(\alpha_{21})} = 0.0000001185757841,\end{aligned}$$

and the sample joint proportions for male are

$$\begin{aligned}
 p_{11}(x_2) &= \frac{\exp(\alpha_{11} + \beta_{11})}{1 + \exp(\alpha_{11} + \beta_{11}) + \exp(\alpha_{12} + \beta_{12}) + \exp(\alpha_{21} + \beta_{21})} = 0.3682333647 \\
 p_{12}(x_2) &= \frac{\exp(\alpha_{12} + \beta_{12})}{1 + \exp(\alpha_{11} + \beta_{11}) + \exp(\alpha_{12} + \beta_{12}) + \exp(\alpha_{21} + \beta_{21})} = 0.1318244456 \\
 p_{21}(x_2) &= \frac{\exp(\alpha_{21} + \beta_{21})}{1 + \exp(\alpha_{11} + \beta_{11}) + \exp(\alpha_{12} + \beta_{12}) + \exp(\alpha_{21} + \beta_{21})} = 0.1321661813 \\
 p_{22}(x_2) &= \frac{1}{1 + \exp(\alpha_{11} + \beta_{11}) + \exp(\alpha_{12} + \beta_{12}) + \exp(\alpha_{21} + \beta_{21})} = 0.3677760084
 \end{aligned}$$

The male did not substantially change their preference due to 319 presidential gun shot, but it affects the final decision of female.

6. SUMMARY

Certain evidence exists to show that the randomized response sampling model reduces bias caused by nonsampling errors which might otherwise result when using direct sampling methods to survey sensitive topics. The theory underlying the application of the randomized response procedure in a matched pair study is presented and applied to data collected for that purpose in a survey of 319 presidential gun shot in Taiwan. Estimates of the cell probabilities and testing the marginal homogeneity are reported in this article. Additionally, the baseline categorical logit model is presented.

References

- [1] Abul-Ela, A.-L., Greenberg, B. G., and Horvitz, D. G. (1967) "A multi-proportions randomized response model" *Journal of the American Statistical Association*, 62, 990-1008.
- [2] Greenberg, B. G., Abul-Ela, A.-L., Simmons, W. R., and Horvitz, D. G. (1969) "The unrelated question randomized response model: Theoretical framework" *Journal of the American Statistical Association*, 64, 520-539.
- [3] Greenberg, B. G., Kuebler, R. R., Jr., Abernathy, J. R., and Horvitz, D. G. (1971) "Application of the randomized response technique in obtaining quantitative data" *Journal of the American Statistical Association*, 66, 243-250.
- [4] Scheers, N. J., and Dayton, C. Mitchell (1988) "Covariate randomized response models" *Journal of the American Statistical Association*, 83, 969-974.
- [5] Warner, S. L. (1965) "Randomized response: A survey technique for eliminating evasive answer bias" *Journal of the American Statistical Association*, 60, 63-69.