*Original Article*

# Least-MSE calibration procedures for corrections of measurement and misclassification errors in generalized linear models

Parnchit Wattanasaruch, Veeranun Pongsapukdee*, and Pairoj Khawsithiwong

*Department of Statistics, Faculty of Science,*
*Silpakorn University, Sanam Chandra Palace Campus, Mueang, Nakhon Pathom, 73000 Thailand.*

**Abstract**

The analyses of clinical and epidemiologic studies are often based on some kind of regression analysis, mainly linear regression and logistic models. These analyses are often affected by the fact that one or more of the predictors are measured with error. The error in the predictors is also known to bias the estimates and hypothesis testing results. One of the procedures frequently used to handle such problem in order to reduce the measurement errors is the method of regression calibration for predicting the continuous covariate. The idea is to predict the true value of error-prone predictor from the observed data, then to use the predicted value for the analyses. In this research we develop four calibration procedures, namely probit, complementary log-log, logit, and logistic calibration procedures for corrections of the measurement error and/or the misclassification error to predict the true values for the misclassification explanatory variables used in generalized linear models. The processes give the predicted true values of a binary explanatory variable using the calibration techniques then use these predicted values to fit the three models such that the probit, the complementary log-log, and the logit models under the binary response. All of which are investigated by considering the mean square error (MSE) in 1,000 simulation studies in each case of the known parameters and conditions. The results show that the proposed working calibration techniques that can perform adequately well are the probit, logistic, and logit calibration procedures. Both the probit calibration procedure and the probit model are superior to the logistic and logit calibrations due to the smallest MSE. Furthermore, the probit model-parameter estimates also improve the effects of the misclassification explanatory variable. Only the complementary log-log model and its calibration technique are appropriate when measurement error is moderate and sample size is high.

**Keywords:** calibration techniques, misclassification, generalized linear models, regression calibration, logistic, logit, probit, complementary calibration procedures

## 1. Introduction

In nonlinear and generalized linear models, the response $Y$ is generally in terms of explanatory variables or predictors $X$ and a covariate $Z$. Such covariate may represent those predictors measured without error for all practical purposes but those for $X$ possibly cannot be exactly observed for all study subjects. In assessing measurement error, atten-

tion needs to be given by a type and a nature of error as well as sources of data which allow modeling of this error. The analyses of model using an unobserved explanatory $X$ can often use only an observable $W$ which is related to $X$, and that $W = X + U$, where $U$ is the measurement error. Thus, model estimators of the response $Y$ in terms of the direct observed predictor $W$ may be poor and biased (Rosner *et al.*, 1989, 1990; Whittermore, 1989; Gustafson and Lee, 2002). Although, a simple and intuitive method such as regression calibration technique is a popular method in measurement error to correct the classical regression model since it is quite easily implemented. However, more complicated calibration

* Corresponding author.
  Email address: veeranun@su.ac.th, veeranun@hotmail.com

techniques in forms of nonlinear models and generalized linear models (GLMs) (Nelder and Wedderburn, 1972; McCullagh and Nelder, 1983; 1989) are still rarely implemented, for example logistic regression that was studied for its measurement error (Rosner *et al.*, 1990; Thoresen and Laake, 2000) and that for its misclassification error (Reade *et al.*, 1991). Consequently, it is important to include measurement error considerations when planning a study, both to enable application of measurement error analysis of data and to ensure validity of conclusions. Moreover, in assessing of model fit and/or its accuracy of parameter estimates for the models through calibration techniques, researchers tend to measure either its bias or the MSE of model parameters, particularly when only one covariate is used in simulation studies (Carroll *et al.*, 1995; Thoresen and Laake, 2000; Gustafson and Nhu, 2002). Other statistics such as deviance statistic for GLMs will be appropriate for model checking, especially when several covariates are included in a model. It is a versatile statistic that is distributed as an asymptotic chi-square. Also its equation form is equivalent to the MSE of estimates and Pearson's chi-square statistic (Thoresen, and Laake, 2000; Agresti, 2002; Lawal, 2003; Ponsapukdee, 2012). For more details see Section 2.3. The basis of a calibration technique is the replacement of *X* by the calibration modeling of the explanatory or the covariate *X* on other related variables, for example (*Z*,*W*) using an approximate working model for the observed data. This procedure seems to be practical and would be very helpful in cases where the observed data are from using replication, validation or instrumental data for the including *X* covariate (Carroll *et al.*, 1995).

In this research, four calibration techniques are developed to predict the true unobserved discrete *X* covariate, $X\_g|W$, from the error-prone observed *W*. Then, the predicted $X\_g$ is used in a case of building the nonlinear models in the form of GLMs to improve the efficiency as well as the accuracy of the GLMs. The main purpose is to investigate the performance of the four proposed calibration techniques, i.e. the probit calibration, the complementary log-log calibration, the logistic calibration, and the logit calibration techniques through GLMs.

The logistic calibration is particularly intended to use with only a continuous explanatory variable $X\_c$. In contrast, the logit calibration is aimed to use only a discrete or a categorical explanatory variable $X\_g$. In fitting GLMs with correcting the measurement error and misclassification error, only three models are desinged, such that the probit, the complementary log-log, and the logit models (Agresti, 2002), because the logit model used here includes only a binary covariate of which the values are predicted from both the logistic calibration and the logit calibration, once at a time. All four calibration procedures and the three model approaches will be investigated considering the deviation of estimates from the true parameters of the models with regard to the mean squared error (MSE) of estimates of the model parameters. A thousand simulations studies at each condition of sample sizes, calibration techniques, model-para-

meters conditions and the GLMs models were performed. All work was processed by using our developed macro program running with the SAS 9.1®.

## 2. Methodology

The article focuses on fitting the statistical models or GLMs relating a binary response *Y*, to data formulated in terms of well-defined but unobservable *X*, using information on measurements *W* that are less than perfectly correlated with *X*. Problems of this nature are called measurement error problems. The statistical models and methods for analyzing such data are called measurement error models. Then, methods are organized in three main issues: 1) models of interest, 2) the proposed calibration techniques, and 3) the assessing of goodness-of-fit.

### 2.1 Models of interest

The fundamental concepts of the three models are defined. These models will be fitted by using the results from the four calibration techniques with the distinct covariates.

#### 2.1.1 Probit model for the binary response

Probit model is defined as

$$\Phi^{-1}\left[P\left(Y=1|x\right)\right]=\beta_0+\beta x, \tag{1}$$

where $\Phi$, is the standard normal cdf. This link function is called the probit link function, $\Phi^{-1}(\cdot)$, and the parameters $\beta_0$ and $\beta$ are model parameters. The model in (l) was fitted using the four calibration techniques separately.

#### 2.1.2 Complementary log-log model for the binary response

Complementary log-log model is given

$$\log\left\{-\log\left[1-P\left(Y=1|x\right)\right]\right\}=\beta_0+\beta x. \tag{2}$$

Then, it can be written in a form,

$$P\left(Y=1|x\right)=1-\exp\left[-\exp\left(\beta_0+\beta x\right)\right].$$

The complementary log-log link is in a form of $\log\left\{-\log\left[1-P\left(Y=1|x\right)\right]\right\}$ since the log-log link applies to the complement of $P\left(Y=1|x\right)$. It is asymmetric and $P\left(Y=1|x\right)$ approaching zero fairly slowly but approaching 1 quite sharply. On the other hand, the logit and probit links are symmetric about 0.5. The model in (2) was fitted using the four calibration techniques separately.

#### 2.1.3 Logit model for the binary response

When an GLMs covariates *X* consists of at least one or all categorical explanatory variables, it is usually called the model as a logit GLM given in (3). By contrast, the name

logistic model which has similar form as (3) usually permits at least one or all continuous explanatory variables in the model (Agresti, 2002).

$$\log \frac{P\left(Y=1|x\right)}{1-P\left(Y=1|x\right)} = \beta_0 + \beta x, \tag{3}$$

where, $\beta_0$ and $\beta$ are model parameters and the link function is called the logit link.

For the model in (3), the logit model is fitted by a binary predicted covariate which is carried out from using both the logit calibration (with a categorical covariate, $w\_g$) and the logistic calibration (with a continuous covariate, $w\_c$), separately. Thus, the model in (3) will also be fitted using the four calibration techniques, separately, once at a time.

The above GLMs differ from the traditional general linear models (for example, a regression model is a special case) in two major aspects. Firstly, the distribution of response variable can be explicitly non-normal, i.e. it can be binomial, Poisson, negative binomial, hypergeometric, multinomial or even product multinomial. Secondly, the response values are predicted from a linear combination of explanatory variables, which are also generalized to mixed categorical and continuous or either of them, and connected to the response variable via a link function. In some situations, the probit link models give the best power of the tests for every test statistic (Pongsapukdee and Sukgumphaphan, 2008). In classical general linear models, the response variable values are expected to follow the normal distribution and the link function is a simple identity function. For GLMs, the response variable follows the exponential family distribution models and the most often used link functions include logit, probit, complementary-log-log, and also the log links.

**2.2 Proposed calibration techniques**

In this part, four calibration procedures was introduced for the correction of the measurement error and the misclassification error to predict the true values for the misclassification explanatory variables used in GLMs as the following.

**2.2.1 Probit calibration procedure**

The probit calibration procedure is somewhat mitigated by the need to develop and to fit a calibration model as the model of $X$ on the other continuous covariate $w\_c$, or the categorical covariate, $w\_g$ as mentioned previously that in practice an unobserved explanatory $X$ often can be observed or collected only an observable $W$ such that $W = X+U$, where $U$ is the measurement error. In the case of probit calibration technique, the observed explanatory $W$ was generated in order to predict or estimate the $X$'s binary values which will be denoted by $x*\_g$. Hence, in predicting the true values of variable $X\_g$ from $w\_c$ or $w\_g$, say $x*\_wc$ or $x*\_wg$, respectively, will be evaluated by using the probit

calibration procedure of the form, for the continuous covariate $w\_c$,

$$\Phi^{-1}\left[P\left(x\_g=1|w\_c\right)\right] = \beta_0 + \beta_1 w\_c$$

and for the categorical covariate ,

$$\Phi^{-1}\left[P\left(x\_g=1|w\_g\right)\right] = \beta_0 + \beta_1 w\_g$$

**2.2.2 Complementary log-log calibration procedure**

Alternatively to those used in the probit calibration technique, by replacing the probit link function to the complementary log-log link, the complementary log-log calibration procedure was obtained by the form, for the continuous covariate $w\_c$,

$$\log\left\{-\log\left[1-P\left(x\_g=1|w\_c\right)\right]\right\} = \beta_0 + \beta_1 w\_c$$

and for the categorical covariate $w\_g$,

$$\log\left\{-\log\left[1-P\left(x\_g=1|w\_g\right)\right]\right\} = \beta_0 + \beta_1 w\_g.$$

**2.2.3 Logistic calibration procedure**

The logistic calibration technique which uses only the continuous explanatory variable $w\_c$, to estimate or to predict the variable $X$, say $x*\_g$, has a form

$$\log\frac{P\left(x\_g=1|w\_c\right)}{1-P\left(x\_g=1|w\_c\right)} = \beta_0 + \beta_1 w\_c.$$

**2.2.4 Logit calibration procedure**

The logit calibration technique which uses only the categorical or grouped explanatory variable $w\_g$, to estimate or to predict the variable $X$, say $x*\_g$, has a form

$$\log\frac{P\left(x\_g=1|w\_g\right)}{1-P\left(x\_g=1|w\_g\right)} = \beta_0 + \beta_1 w\_g.$$

**2.3 Assessing of goodness of fit**

In the context of GLMs, likelihood ratio model comparison using the deviance is usually investigated by considering two models: $M_0$ with fitted values $\hat{\mu}_0$, and $M_1$ with $\hat{\mu}_1$, with $M_0$ a special case which is nested within $M_1$. The likelihood ratio test ($G^2$) (or also the log-likelihood ratio test) was originally defined by Wilks in 1938, where $G^2 =$

$$2\sum \text{observed} \times \log\left(\frac{\text{observed}}{\text{fitted}}\right) = 2\left[\log L\left(\hat{\mu}_0;y\right) - \log L\left(\hat{\mu}_1;y\right)\right].$$

The deviance ($D$) was also originally introduced by Nelder and Wedderburn (1972), where $D = -2\log\left[\log L\left(\hat{\mu}_0;y\right) - \right.$ $\left. \log L\left(\hat{\mu}_1;y\right)\right]$. Thus, the deviance is equivalent to the likelihood ratio test by definitions and the simpler models have larger deviances. A model based on $p$ parameters with $n$

observations would have its computed test statistic distributed as $\chi^2_{n-p}$. For example, in a case of Poisson observations $D$ can be directly calculated. However, for some other distributions, $D$ may be indirectly computed directly in spite of any nuisance parameters. For the normal distribution, it can be shown that $\sigma^2 D = \sum (y_i - \hat{\mu}_i)^2$, where $\hat{\mu}_i$ denotes the MLE of $\mu_i$. The PROC GENMOD in SAS$^{®}$ can obtain $\sigma^2 D = \sum (y_i - \hat{\mu}_i)^2$ and gives the scale parameter which is an estimate of $\sigma^2$ in term of MSE, i.e. $\hat{\sigma}^2 = \dfrac{D}{n-p}$. Therefore, the term MSE can be obtained by the deviance $D$. For the one way multinomial, both $G^2$ and $D$ have an asymptotic $\chi^2$ distribution with $p-1$ degrees of freedom in the case of specified probabilities. In this simulation study under model conditions with one covariate and known model parameter, the MSE of the estimated model parameter is straight forward computed and investigated for 1,000 sets by comparing among the least-MSE of all combinations of calibration techniques and models of interest in each condition.

## 3. Simulation Experiments

From the models and the calibration techniques in Section 2, the simulation studies were conducted for the dichotomous response categories Y with the model parameters: $\beta_0 = -2.25$ and $\beta_1 = 0.371$ (Thoresen and Laake, 2000). For the measurement error terms, $U$ is generated from $N\left(0, \sigma_U^2\right)$ where $\sigma_U^2 = 0.75$, 1, and 3. The explanatory variables $X\_c$ is from $N(0,1)$. The continuous observable covariate $W\_c$ is in a term of $W\_c = X\_c + U$ and the categorical observable covariate $W\_g$ is obtained by if $W\_c > 0$, then $W\_g = 1$ and $W\_c = 0$ elsewhere. Data are simulated under the sample sizes of 100, 500, and 1,000, according to that the samples needed

to achieve the power 0.90-0.95, and when using the Bernoulli (0.5) explanatory variable, the sample units would be closing 1,000 (Shieh, 2001). In each condition of sample sizes and that of parameters of covariates' distributions and the model parameters, the calibration techniques are performed to obtain the probabilities $P(x\_c)$ and the probabilities $P(x\_g)$. Then, the $X^*\_(\cdot)$ values can be attained by $x^*\_wc = 1$ if $P(x\_c) > u$ and $x^*\_wc = 0$ elsewhere. As a same fashion, $x^*\_wg = 1$ if $P(x\_g) > u$ and $x^*\_wg = 0$ elsewhere, where $u$ is from $U(0,1)$. Therefore, the response outcomes $y_j, j = 1,2,...,n$ and the $P(Y = 1|x)$ estimates were then computed through the correctly specified models in each case of sample size. Each condition was carried out for 1,000 repeated simulations using random response outcomes of $Y$ with the same set of all $X^*\_(\cdot)$ values from the calibration techniques. The developed macro was run on SAS 9.1$^{®}$. Statistical analyses for assessing the accuracy of the models and all calibration techniques are based on the MSE statistics obtained from the models fitted under their corresponding conditions.

## 4. Results

The results in terms of the minimum value of MSE from the comparison among the four calibration techniques are evaluated for each condition and model fitted. The calibration procedures with $w\_c$ or $w\_g$ under $\sigma_U^2 = 0.75$, classified by the sample size and the model fitted indicated that the logistic procedure can give the least MSE (0.07389) under the probit model with the continuous covariate ($w\_g$) for the sample size of 100 (Table 1). Meanwhile, the logit calibration procedure provides the least MSE (0.01063, 0.00592) with the categorical covariate ($w\_c$) for the sample sizes of 500 and 1,000, respectively (Table 1). Therefore, when the $\sigma_U^2 = 0.75$ the minimum MSE is from the logit calibration procedure under the probit model (Table 1 and Figure 1).

Table 1. Least-MSE calibration procedure with $w\_c$ or $w\_g$ under $\sigma_U^2 = 0.75$ classified by the sample size and the model fitted.

| $\sigma_U^2$ | Sample size | Model | Calibration Procedure | |
|---|---|---|---|---|
| | | | $w\_c$ | $w\_g$ |
| 0.75 | 100 | Probit | **Logistic (0.07389)** | Comp-log-log (0.08029) |
| | | Comp-log-log | Logistic (0.10414) | Comp-log-log (0.11015) |
| | | Logit | Logistic (0.18961) | Comp-log-log (0.20601) |
| | 500 | Probit | Logistic (0.01084) | **Logit (0.01063)** |
| | | Comp-log-log | Logistic (0.01363) | Logit (0.01338) |
| | | Logit | Logistic (0.02769) | Logit (0.02717) |
| | 1,000 | Probit | Logistic (0.00607) | **Logit (0.00592)** |
| | | Comp-log-log | Logistic (0.00772) | Logit (0.00752) |
| | | Logit | Logistic (0.01548) | Logit (0.01511) |

Each value in parenthesis is the least MSE value, or the minimum value which is compared among four calibration techniques under the same set of $\sigma_U^2$, Sample size and Model.
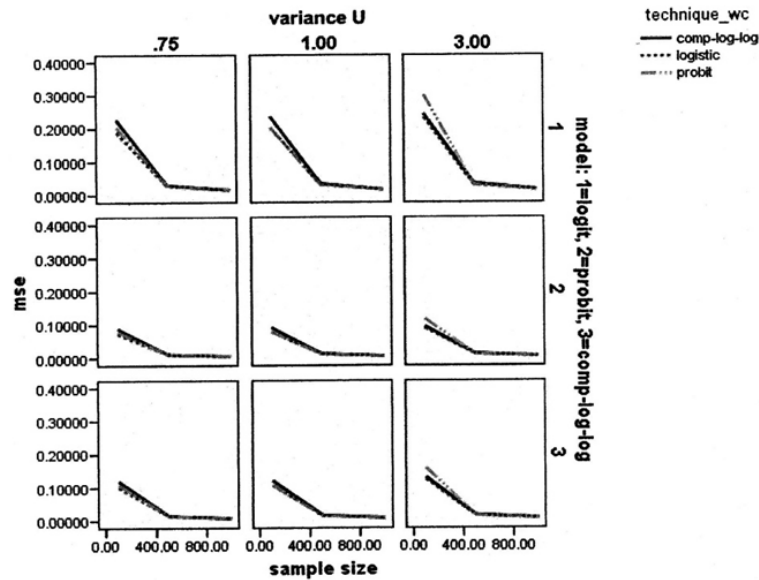
Figure 1. MSE plots of calibration techniques with the continuous covariate $w\_c$, classified by the sample size, variance $U\left(\sigma_U^2\right)$ and the model fitted (model).

Similarly, comparison results when $\sigma_U^2 = 1.00$ and $\sigma_U^2 = 3.00$ showed that the minimum MSE is from the probit calibration procedure under the probit model, for both the sample size of 500 and 1000 with MSE = 0.01096, 0.00607 (Table 2) and MSE = 0.01137, 0.00619 (Table 3), respectively. However, the next smallest MSE is from the logistic calibration procedure under the probit model, MSE = 0.09056 (Table 3). Therefore, for the final results, it is shown that the proposed probit calibration procedure is probably chosen and would be the most appropriate procedure for the generalized linear models under the probit model when the analysis considering the measurement error and misclassification error (Table 2-3 and Figure 2).

## 5. An application of the propose procedures

In this section, the results of the proposed tests on the calibrations techniques and the models of interest are presented to show their application in a real example on antibiotics/SIDs data from Greenland (1988). The data involve a case-control study of the association of antibiotic use by mother during pregnancy $X$, and the occurrence of sudden infant death syndrome (SIDs) $Y$. The error-prone measurement of antibiotic use ($W$) is based on a self report from mother. The main study data of 428 women are randomly selected to be a training set. Another validation set of 428 women are also determined the true antibiotic use ($X$) for

Table 2. Least-MSE calibration procedure with $w\_c$ or $w\_g$ under $\sigma_U^2 = 1.00$ classified by the sample size and the model fitted.

| $\sigma_U^2$ | Sample size | Model | Calibration Procedure | |
| --- | --- | --- | --- | --- |
| | | | $w\_c$ | $w\_g$ |
| 1.00 | 100 | Probit | Logistic (0.07944) | **Comp-log-log (0.07604)** |
| | | Comp-log-log | Logistic (0.10687) | Comp-log-log (0.10386) |
| | | Logit | Logistic (0.20344) | Comp-log-log (0.19478) |
| | 500 | Probit | Probit (0.01187) | **Probit (0.01096)** |
| | | Comp-log-log | Logistic (0.01493) | Probit (0.01383) |
| | | Logit | Probit (0.03033) | Probit (0.02799) |
| | 1,000 | Probit | Comp.log-log (0.00611) | **Probit (0.00607)** |
| | | Comp-log-log | Comp.log-log (0.00776) | Probit (0.00811) |
| | | Logit | Comp.log-log (0.01559) | Probit (0.01550) |

Each value in parenthesis is the least MSE value, or the minimum value which is compared among four calibration techniques under the same set of $\sigma_U^2$, Sample size and Model.

Table 3.   Least-MSE calibration procedure with $w\_c$ or $w\_g$ under $\sigma_U^2 = 3.00$ classified by the sample size and the model fitted.

| $\sigma_U^2$ | Sample size | Model | Calibration Procedure | |
| --- | --- | --- | --- | --- |
| | | | $w\_c$ | $w\_g$ |
| 3.00 | 100 | Probit | **Logistic (0.09056)** | Probit (0.09440) |
| | | Comp-log-log | Comp-log-log (0.13113) | Probit (0.12944) |
| | | Logit | Logistic (0.23261) | Probit (0.24185) |
| | 500 | Probit | Logistic (0.01149) | **Probit (0.01137)** |
| | | Comp-log-log | Logistic (0.01448) | Probit (0.01431) |
| | | Logit | Logistic (0.02936) | Probit (0.02905) |
| | 1,000 | Probit | Logistic (0.00624) | **Probit (0.00619)** |
| | | Comp-log-log | Logistic (0.00792) | Probit (0.00787) |
| | | Logit | Logistic (0.01591) | Probit (0.01579) |

Each value in parenthesis is the least MSE value, or the minimum value which is compared among four calibration techniques under the same set of $\sigma_U^2$, Sample size and Model.
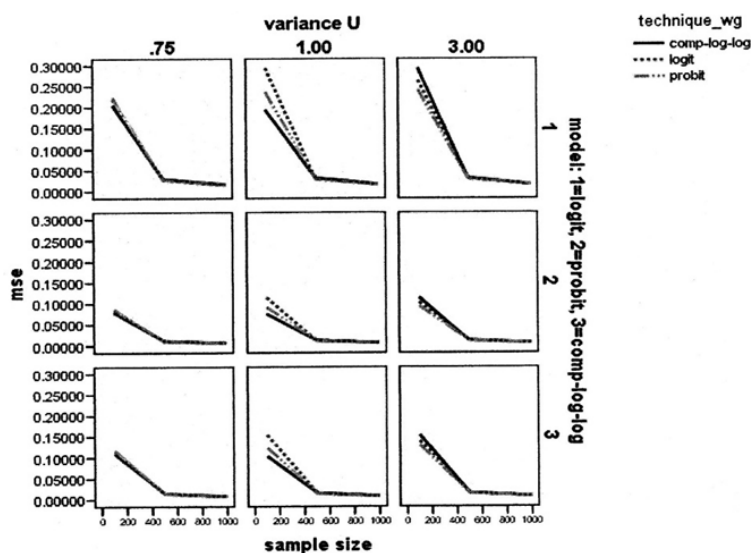


Figure 2.   MSE plots of calibration techniques with the continuous covariate $w\_g$, classified by the sample size, variance $U\left(\sigma_U^2\right)$ and the model fitted (model).

women as determined from medical records (Table 4). The analysis results from the application of the proposed calibrations procedures on this real clinical trial example provides an example in practice that the proposed probit calibration technique demonstrates the best properties in terms of the estimated MSE and deviance (Table 5). This is consistent with the rational underlying the proposal and it is confirmed by the simulation studies in Section 3.

## 6. Conclusions

In conclusion, the results show that the calibration techniques that perform adequately well, are respectively from the probit calibration procedure with both $w\_c$ and $w\_g$ covariates, the logistic procedure with $w\_c$, and the logit calibration procedure, with $w\_g$. The probit calibration procedure and the probit model is superior to the logistic and logit calibration procedures due to the smallest MSE. Furthermore, the probit model parameter estimates do improve the effects of the misclassification explanatory variable. Hence, the results indicate that use of the calibration procedures in generalized linear models, the probit calibration procedure has most statistically accuracy results and is probably able to use safely. In addition, it is shown that the logistic calibration procedure is appropriate with $w\_c$ and the logit calibration procedure is appropriate with $w\_g$.

Table 4.  Application of real data with training and validation data sets.

| | | Control (Y = 0) X | | Cases (Y = 1) X | |
|---|---|---|---|---|---|
| | | 0 (no use) | 1 (use) | 0 (no use) | 1 (use) |
| **Training Data Set** | | | | | |
| W | 0 (no use) | 94 | 88 | 76 | 95 |
| | 1 (use) | 22 | 17 | 16 | 20 |
| | | 116 | 105 | 92 | 115 |
| **Validation Data Set** | | | | | |
| W | 0 (no use) | 168 | 16 | 143 | 17 |
| | 1 (use) | 12 | 21 | 22 | 29 |
| | | 180 | 37 | 165 | 46 |

Table 5.  Application results of estimated least MSE and deviance under the best calibration procedure classified by the model fitted.

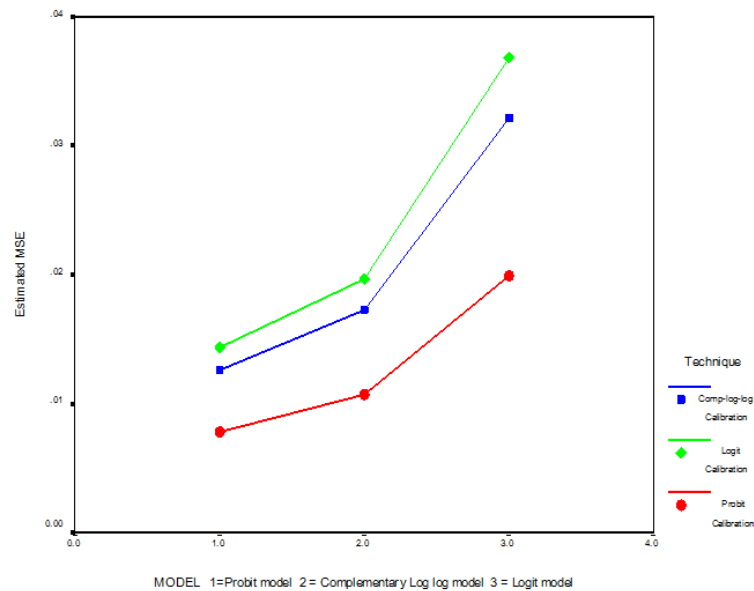| Model | Calibration Procedure | Deviance |
|---|---|---|
| Probit | Probit Calibration (0.0078) | 591.993 |
| Logit | Probit Calibration (0.0144) | 592.417 |
| Comp-log-log | Probit Calibration (0.0126) | 592.328 |



Figue 3.  The estimated MSE plots of the three calibration techniques from the discrete real data set classified by the model fitted.
Note : There are only three calibration techniques since all real data set used in this application are discrete or categorical data.

However, only the complementary log-log model and its calibration procedure are appropriate when the measurement error is moderate and the sample size is large.

## References

Agresti, A. 2002. Categorical Data Analysis. 2nd edition, John Wiley & Sons. New York, U.S.A.

Buonaccorsi, J.P. 2010. Measurement Error: Models, Methods, and Applications. Chapman & Hall. New York, U.S.A.

Carroll, R.J., Ruppert. D. and Stefanski, L.A. 1995. Measurement Error in Nonlinear Models. Chapman & Hall. London, U.K.

Greenland, S. 1988. Variance estimation for epidemiologic effect estimates under mis-classification. Statistics in Medicine. 7, 745-758.

Gustafson, P, and Nhu D.Le. 2002. Comparing the effects of continuous and discrete covariate mismeasurement, with emphasis on the dichotomization of mismeasured predictors. Biometrics. 58, 878-887.

Lawal, H.B. 2003. Categorical Data Analysis with SAS and SPSS Applications. Lawrence Erlbaum Associates. Inc. London, U.K.

McCullagh, P. and Nelder, J.A. 1983. Generalized Linear Models. 2nd edition1989, Chapman & Hall, London, U.K.

Nelder, J.A. and Wedderburn, R.W.M. 1972. Generalized linear models. Journal Royal Statistics A. 135, 370-384.

Pongsapukdee, V. and Sukgumphaphan, S. 2008. Three cumulative link models for ordinal responses: A review of methods and power comparisons. Songklanakarin Journal of Science and Technology. 30(6), 805-811.

Pongsapukdee, V. 2012. Analysis of Categorical Data: Theories and Applications with GLIM, SPSS, SAS and MTB. 3rd edition, Silapakorn University Press. Nakhon Pathom, Thailand.

Rosner, B., Spiegelman, D. and Willett, W. C. 1990. Correction of logistic regression relative risk estimates and confidence intervals for measurement error: The case of multiple covariates measured with error. American Journal of Epidemiology. 132, 734-745.

Reade, C., Susan. J. and Kupper, L. 1991. Effects of exposure misclassification on regression analysis of epidemiologic follow-up study data. Biometric. 47 (2), 535-548.

Rosner, B., Willett, W. C. and Spiegelman, D. 1989. Correction of logistic regression relative risk estimates and confidence intervals for systematic within-person measurement error. Statistics in Medicine. 8, 1051-1070.

Shieh, G. 2001. Sample size calculations for logistic and Poisson regression models. Biometrika. 88 (4), 1193-1199.

Thoresen, M. and Laake, P. 2000. A simulation study of measurement error correction methods in logistic regression. Biometrics. 56, 868-872.

Whittemore, A.S. 1989. Errors in variables regression using Stein estimates. American Statistician. 43, 226-228.