

## GENERALIZED VARIANCE FUNCTIONS FOR BINOMIAL VARIABLES IN STRATIFIED TWO-STAGE SAMPLING

Ari Handayani<sup>1)</sup>, Aunuddin<sup>2)</sup> and Indahwati<sup>2)</sup>

<sup>1)</sup> Directorate of Statistical Methodology, BPS-Statistics Indonesia

<sup>2)</sup> Department of Statistics, FMIPA IPB

### Abstract

*This empirical study evaluates the application of Generalized Variance Functions (GVFs) for binomial variables in the 1998 Indonesian Labor Force Survey. The survey employs stratified two-stage cluster sampling for selecting samples from a population of households. The study covers all provinces in Java to produce estimates at the level of Java Island. The relative variance estimates resulted from the GVF models are compared to the relative variance estimates which are computed directly.*

*The results illustrate that model  $\hat{u} = c\hat{Y}^d$  expressed by logarithmic model  $\log \hat{u} = \log c + d \log (\hat{Y})$  gives a good approximation to estimate the variances for the nonagricultural employment group, especially for working male category both in urban and rural areas. It is also good for the total employment group differentiated by age group, educational attainment, and employment status. On the other hand, the model gives poor results for the agricultural employment group.*

*Based on the empirical results, the GVF models may not perform particularly well for the common characteristics which have relatively dissimilar deff values to majority of characteristics in the same group, since these characteristics usually come out among all persons in the sample household and often among all households in the sample cluster as well. The success of the GVF technique depends critically on the grouping of the estimates total ( $\hat{Y}$ ) and amount of characteristics involved as the observations for fitting the model. Furthermore, observations with relatively large residuals will also determine the performance of goodness-of-fit of the model.*

*Application of GVF technique to obtain an approximate standard error on numerous binomial characteristics in large scale survey should be carried out further using extensive data. The better performance of GVF model may also be accomplished by utilizing, for examples, weighted least squares procedure or robust regression method. Additionally, the data users should be warned that there will inevitably be survey characteristics for which GVF's will give poor results or even no GVF will be appropriate.*

Keywords : Generalized Variance Functions, Stratified Two-Stage Sampling

### INTRODUCTION

The measurement of precision for all survey estimates has become a basic and principal need in almost all survey analyses. Variance estimate is mostly used to measure this precision. Variance estimates in large scale survey, where statistics are published for many characteristics, for each of several demographic subgroups of the total population and possibly for a number of geographic areas, can be produced for all survey characteristics simply by evaluating a model at the

survey estimates, rather than by direct computation.

There are practical reasons why general techniques are more desirable. Presentation of individual sampling error would usually more costly and time consuming. Besides it would essentially double the size of tabular publication. These considerations have led to the use of model as a means of approximating sampling error. This model expresses the variance as a function of the expected value of the survey estimate. This method of variance estimation is called the

method of Generalized Variance Function (GVF) (Wolter, 1985).

This study intends to derive the GVF models for several employment groups in the 1998 Indonesian Labor Force Survey. The study also compares the relative variance estimates which are resulted from the GVF models to the relative variance estimates which are computed directly.

This paper test the theory of GVF empirically. Section Theory presents sample design and estimation methods in the 1998 Indonesian Labor Force Survey. This section also describes a class of models for which a GVF estimator is appropriate. Section Methodology describes source of data and analytic method used in this empirical study. Section Results and Discussion summarizes results of the empirical study of GVF as applied to estimator of the total of binomial variables. The last section concludes the paper with a brief summary.

## THEORY

### Sample Design and Estimation Methods

The 1998 Indonesian Labor Force Survey is especially designed to look at a shift of labor force structure between enumeration periods, besides to provide information on manpower statistics and key indicators of labor market at national level. Its results are disseminated in the issue of *Labor Force Situation in Indonesia, August 1998* (BPS-Statistics Indonesia, 1998b) which contains about 300 pages of tabular publication.

The 1998 Indonesian Labor Force Survey employs stratified two-stage cluster sampling for selecting samples from a population of households. The enumeration areas (EAs) are stratified into urban and rural. At the first stage, a sample of  $n_h$  EAs is selected from stratum  $h$ , using linear systematic sampling. At the second stage a subsample of one segment group is selected from the total of segment groups in the selected EAs, by probability proportional to size sampling. At the third stage, a subsample of 12 households is selected from the list frame of households in the selected segment groups, using linear systematic sampling. Since only one segment group is selected for each selected EA at the second stage, then the number of selected EAs are same as the number of selected segment groups. This stage is called the effective sampling stage. Therefore, the sample design is regarded as stratified two-stage cluster sampling, where each segment group, that used as a cluster, is considered as a primary sampling unit, and each household is associated with an ultimate sampling unit, with urban and rural areas are considered as the strata.

Estimate for total values of each characteristic for stratum of urban or rural area at the provincial level is calculated using indirect estimation, *ratio estimate*,

$$\hat{Y}_h = \frac{\hat{P}_h}{12n_h} \sum_{i=1}^{n_h} \sum_{j=1}^{12} \frac{1}{A_{hij}} \sum_{k=1}^{A_{hij}} y_{hijk} = \hat{P}_h \bar{y}_h \quad (1)$$

where,

$\hat{Y}_h$  = estimate of total characteristic  $y$  for stratum  $h$ ,

$\hat{P}_h$  = estimate for population number as a result of population projection (aged 15 years and over) for stratum  $h$ ,

$\bar{y}_h$  = sample means value of characteristic  $y$  for stratum  $h$ ,

$A_{hij}$  = number of household members at household  $j$ , cluster  $i$ , for stratum  $h$ ,

$y_{hijk}$  = characteristic value  $y$  of household member  $k$  (aged 15 years and over), at household  $j$ , cluster  $i$ , for stratum  $h$ .

Estimate for total values of each characteristic at the provincial level is

$$\hat{Y} = \hat{Y}_U + \hat{Y}_R \quad (2)$$

where,

$\hat{Y}$  = estimate of total characteristic  $y$ ,

$\hat{Y}_U$  = estimate of total characteristic  $y$  for urban area,

$\hat{Y}_R$  = estimate of total characteristic  $y$  for rural area.

The variance estimate consider as

$$s_{\hat{Y}}^2 = \sum_{h=1}^H \left[ \frac{n_h}{n_h - 1} \sum_{i=1}^{n_h} \left[ \hat{Y}_{hi} - \frac{\hat{Y}_h}{n_h} \right]^2 \right] \quad (3)$$

where  $\hat{Y}_{hi} = \hat{P}_{hi} \bar{y}_{hi}$  and  $n_h$  is number of selected clusters in stratum  $h$ .

Let  $\hat{Y}$  denote an estimator of total characteristic  $y$  and let  $Y = E(\hat{Y})$  denote its expectation. Consider  $\sigma_{\hat{Y}}^2 = Var(\hat{Y})$  denote the variance of  $\hat{Y}$ , then  $u = Var(\hat{Y})/Y^2$  is called the relative variance of  $\hat{Y}$  (Wolter, 1985).

It is useful to compute the ratio of the two variance estimates, the proper estimate divided by the estimate of a simple random sample of the same sample size,  $n$ . This ratio measures the design effect or *deff* (Kish, 1995).

$$deff = s_{\hat{Y}}^2 / s_{srs}^2 \quad (4)$$

### Generalized Variance Functions

Generalized Variance Function is a simple model that expresses variance as a function of the expected value of the survey estimate. The GVF that is used to estimate the variance of an estimated population total,  $\hat{Y}$ , is of the form

$$\text{Var}(\hat{Y}) = aY^2 + bY \quad (5)$$

The model is formed based on the assumption that the variance of  $\hat{Y}$  can be expressed as product of variance from simple random sample for a binomial random variable and a design effect (U.S. Bureau of Labor Statistics, 1996a).

Consider  $P=Y/N$  as the proportion of the population having characteristic  $Y$ , where  $N$  is the population size and  $Q=1-P$ . The variance of the estimated total  $\hat{Y}$  based on a sample of  $n$  individuals from the population, is

$$\text{Var}(\hat{Y}) = \frac{N^2 PQ(\text{deff})}{n} \quad (6)$$

This can be written as

$$\text{Var}(\hat{Y}) = -(\text{deff}) \left( \frac{N}{n} \right) \left( \frac{Y^2}{N} \right) + (\text{deff}) \left( \frac{N}{n} \right) Y$$

If  $a = -\frac{b}{N}$  and  $b = \frac{(\text{deff})N}{n}$  then

$$\text{Var}(\hat{Y}) = aY^2 + bY$$

Defining  $u = \text{Var}(\hat{Y})/Y^2$  as the relative variance, then model (5) can be written as

$$u = a + \frac{b}{Y} \quad (7)$$

The parameters  $a$  dan  $b$  are estimated by fitting the model to a group of related estimates of total ( $\hat{Y}$ ) and their estimated relative variances ( $\hat{u}$ ).

Valliant (1987) has introduced an alternative for the relative variance in (7) which can be stated as

$$u = -\frac{(\text{deff})}{n} + \frac{(\text{deff})N}{nY}$$

$$u = \frac{(\text{deff})N}{nY} \left[ 1 - \frac{Y}{N} \right] = \frac{(\text{deff})N}{nY} [1 - P]$$

If the proportion of the population having characteristic  $Y$ ,  $P$  is small then this approximation suggests an alternative model  $u = cY^d$  which can be expressed as logarithmic model

$$\log u = \log c + d \log (Y) \quad (8)$$

The success of the GVF technique depends critically on the grouping of the survey statistics, whether all statistics within a group behave according to the same mathematical model or not. This implies that all statistics within group should have a common design effect,  $\text{deff}$ . From a substantive point of view, the grouping will often

be successful when the statistics, (1) refer to the same basic demographic or economic characteristic, (2) refer to the same race-ethnicity group, and (3) refer to the same level of geography (Wolter, 1985). This should give us estimates in the same group that have similar design effects.

The final groups then can be evaluated using scatterplot of  $\hat{Y}$  versus  $\hat{u}$ . Next, the sample coefficient of determination,  $R^2$ , is used to determine the proper model. When all data in the scatterplot lie on the least square line,  $R^2 = 1$  (Smith, 1998). The observed residuals are also key indicators of goodness-of-fit of the model. If the plotted points of residuals against the fits  $\hat{Y}$  are randomly scattered about the horizontal line  $e=0$ , then it explain that the simple linear regression model is appropriate. However, if the residual plot is trumpet-shaped, then nonconstant residual variance is diagnosed. If the residual plot exhibits a curved pattern, it indicates that the analysis based on simple linear regression is incorrect (Aunuddin 1989, Weisberg 1985).

### The Application of GVFs

Two major surveys in the United States, Current Population Survey (CPS) and National Health Interview Survey (HIS) use GVF model  $u = a + \frac{b}{Y}$  (Valliant 1987, Wolter 1985). This model has been used in the CPS since 1947 and is applied to cases in which  $Y$  is the total units that have some binomial characteristics. In those surveys, the  $a$  and  $b$  parameters are estimated using an iterative reweighted least squares procedure, the weight is  $1/u^2$ . Periodically the  $a$  and  $b$  parameters are updated to reflect changes in the ratio between population size and sample size,  $N/n$ . This can be done without recomputing direct estimate of variances as long as the sample design and estimation procedures are essentially unchanged (U.S. Bureau of Labor Statistics, 1996a).

The GVF methods are mainly applicable to the problem of variance estimation for an estimated proportion or for an estimate of the total number of individuals. There have been a few attempts, not entirely successful, to develop GVF techniques for quantitative characteristics. It is very difficult for the quantitative characteristics to have all statistics within a group that behave according to the same mathematical model (Wolter, 1985).

The results in a simulation study, done by Valliant (1987), using household data collected in the CPS show that there are two limitations of using GVF techniques. The limitations are that

they may not perform particularly well for rare characteristics and that there will inevitably be survey variables for which no GVF will be appropriate.

It is also important to keep in mind that standard errors computed from this method reflect contributions from sampling error and some kinds of nonsampling error, and indicate the general magnitude of an estimated standard error rather than its precise value (U.S. Bureau of Labor Statistics, 1996b).

## METHODOLOGY

This simulation study is using household data collected in the 1998 Indonesian Labor Force Survey. The study only considers the working respondents aged 20 - 49 years in Java island which is consisted of four provinces.

The GVFs are applied to the variance estimates of the binomial variables for the labor force status and other demographic characteristics. The binomial variables are formed based on age group (with interval length of 5 years), sex, educational attainment, main industry, kind of occupation, and employment status. Variance estimates for those characteristics are computed for estimation at the level of Java island (not at the provincial level), since the scope of data and characteristics examined are limited. Based on the theoretical justification previously described, the GVF model  $\log u = \log c + d \log (Y)$  is used. The model fitting technique is the least squares method.

### Analytic Method

The simulation study is completed in the following steps.

1. Grouping and selecting the estimates of total characteristic  $\hat{Y}$ .
  - a. Group together all estimates of total characteristic  $\hat{Y}$  that follow a common model, that is based on similar item of characteristics and common deff values.
  - b. Select several members of the group of the estimates of total characteristic  $\hat{Y}$  formed in step a.
  - c. Calculate the estimate of relative variance  $\hat{u} = s_Y^2 / \hat{Y}^2$  for each of the characteristics selected in step b, using direct computation.
  - d. Evaluate the grouping and selection process using scatterplot of  $\hat{Y}$  versus  $\hat{u}$ .
2. Estimating the Parameters of the GVF Models.

Using data  $(\hat{Y}, \hat{u})$ , calculate the estimates of parameters, in relative variance model  $\log \hat{u} = \log c + d \log (\hat{Y})$  using the least squares method.

### 3. Evaluating the GVF Models.

The models resulted are evaluated by diagnostic analyses of the residual plots.

### 4. Estimating the Relative Variance and Standard Error from the GVF models.

The standard error of estimates of  $\hat{Y}$ , which are not computed directly, can be obtained by evaluating the following formula at the survey estimates  $\hat{Y}$ .

$$s_Y^r = \sqrt{c \hat{Y}^{2+d}}$$

### 5. Comparing the Relative Variance Estimates.

The relative variance estimates resulted from the GVF models are compared to the relative variance estimates computed directly. The comparison is carried out by checking the consistency of the ratios between the two relative variance estimates which are expected around the horizontal line  $r = 1$ .

The deff values and relative variance estimates, which are directly calculated, are computed using *IMPS* package program version 3.1 (U.S. Bureau of the Census, 1995). The regression analyses are obtained using *minitab* 11 package program (Minitab Inc, 1996).

## RESULTS AND DISCUSSION

### Grouping and Selecting the Characteristics Prior to Model Estimation

For each of 34,854 working respondents aged 20 - 49 years in Java island, data on labor force status and other demographic characteristics are coded into 36 binomial variables. A respondent's value of a binomial variable is 1 if the respondent has a particular characteristic, and is 0 if not. The binomial variables are formed based on age group, sex, educational attainment, main industry, kind of occupation, and employment status.

In the next discussion, the analyses are distinguished by four demographic categories because statistics in different demographic categories tend to differ with regard to the specific distribution they have. The four demographic categories are :

1. working male in urban area,
2. working male in rural area,
3. working female in urban area, and
4. working female in rural area.

Regarding the similarity of characteristics and common deff values, all estimates of total

characteristic  $\hat{Y}$  are divided into three groups with model (8) fitted independently in each group for each demographic category. Thus, different estimated parameters are obtained for each of the three groups for each of four demographic categories. The three groups are :

1. Agriculture employment, differentiated by educational attainment and kind of occupation as the agricultural workers.
2. Nonagriculture employment, differentiated by educational attainment and kind of occupation as the nonagricultural workers.
3. Total employment, differentiated by age group, educational attainment, employment status and kind of occupation as the nonagricultural workers.

Separation of agricultural employment group and nonagricultural employment group is made since there is a geographic distribution difference of persons employed in agriculture and persons employed in nonagricultural industries. A separate total employment group is used because statistics in this group tend to have similar design effects.

Main industries observed in the nonagricultural employment group consist of manufacturing industry, trade, hotel and restaurant, transportation, construction and public services. There is no sample of respondent employed in transportation and construction sector for the working female category. Other main industries are not analyzed as the observed characteristics are less. The same matter also prevail for the educational level characteristics especially for Vocational Junior High School and Diploma I/II, kind of occupation characteristic such as manager, and employment status characteristic such as employer.

Agricultural employment group and the agricultural workers tend to have fairly large design effects. The reason for this is that the sample respondent mainly employed in agriculture (35%) as an agricultural worker, hence these characteristics usually appear among all persons in the sample household. On the other side, labor force characteristics differentiated by age group, educational attainment, employment status, and nonagricultural worker have lower design effects, since these characteristics tend to vary among members of the same household and among households within a cluster.

Subsequently, several members of the groups of the estimates of total characteristic  $\hat{Y}$  are selected for which variances are estimated directly. The selections include certain keys of labor force statistics and the need to obtain well-fitting GVs that pertain to all statistics. For the agricultural employment group, the observed

characteristics vary from 22 into 32 characteristics since the available characteristics are limited for each demographic category. In the nonagricultural employment group, 75 and 50 characteristics are selected for the working male and working female categories, respectively. Lastly, there are 100 characteristics selected for each demographic category for the total employment group. Then the relative variance estimates  $\hat{u}$  are computed directly for each of the characteristics selected in the former step.

### Fitting the Model

The earlier grouping and selection process are then evaluated using scatterplot of the estimates of total characteristic ( $\hat{Y}$ ) versus its relative variance estimates ( $\hat{u}$ ). Generally, these scatterplots do not indicate a linear model. Suitable transformation of data ( $\hat{Y}, \hat{u}$ ) can sometimes be found that will permit a nonlinear model to be approximated by a linear one. In this case, taking logarithm will reduce this model to a linear model. Then a plot of the log of the estimated relative variance versus the log of the estimate is drawn. Figure 1 shows a log-log plot for working male in the agricultural employment group. The figure is also a plot of the estimated relative variance versus the estimate. The figure displays that some of the observations do not follow the same model as the rest of the data. It comes about given that the characteristics within this group are likely to have deff values that vary considerably from 1.35 into 7.32. Removing the observations that appear to follow a different model cannot be done since the available characteristics are limited for each demographic category.

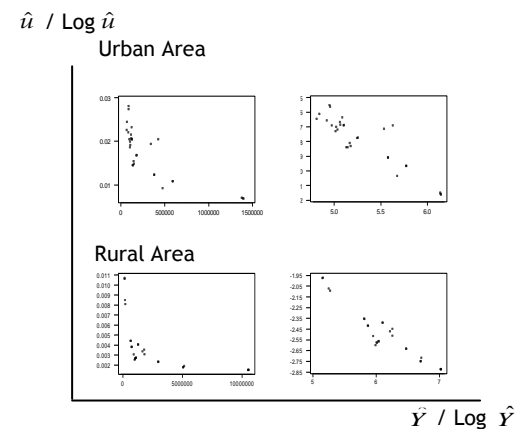


Figure 1. Log-Log Plots of Estimated Relative Variances Versus Estimates for Working Male Category in Agricultural Employment Group

The transformation yields a better linear model for the nonagricultural employment group, except for the working female category in urban and rural areas (see Figure 2). Most of the cases are caused by the characteristics that own relatively large relative variance estimates.

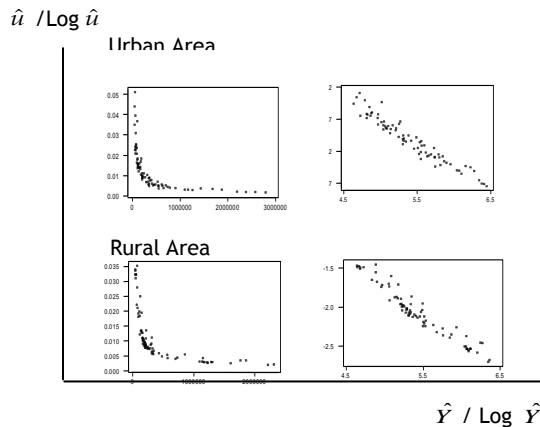


Figure 2. Log-Log Plots of Estimated Relative Variances Versus Estimates for Working Male Category in Nonagricultural Employment Group

Figure 3 shows that the total employment group tends to produce a better linear plot, since the estimates for the characteristics in this group have a tendency to behave according to the same mathematical model as they have similar deff values.

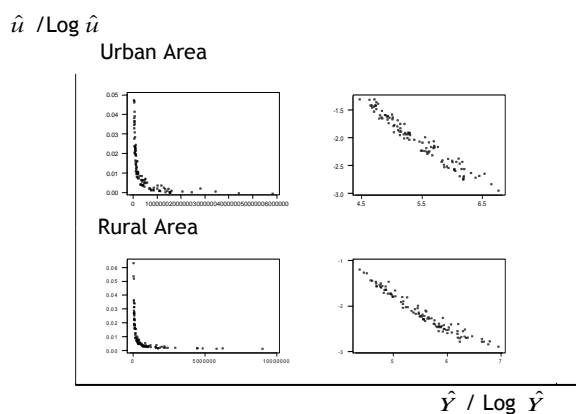


Figure 3. Log-Log Plots of Estimated Relative Variances Versus Estimates for Working Male Category in Total Employment Group

Since the proportions of the population having characteristics observed (P) are generally small, model  $\log \hat{u} = \log c + d \log (\hat{Y})$  is applied for fitting the model. The model parameters are estimated using least squares method. The resulted GVF models are summarized in Table 1.

The results show that the total employment group has the best overall empirical performance. This is indicated by the highest value of the adjusted  $R^2$ .

Table 1. The GVF Models for the Three Groups of Characteristics Distinguished by Four Demographic Categories

Groups of Characteristics Distinguished by Four Demographic Categories	Model $\log \hat{u} = \log c + d \log (\hat{Y})$	Adj. $R^2$	Model $\hat{u} = c\hat{Y}^d$	
			c	d
Agricultural Employment				
Male in Urban Area	$\log \hat{u} = 0.266 - 0.388$	79.8	1.9450	-0.388
Male in Rural Area	$\log \hat{u} = 0.122 - 0.426$	88.9	1.3743	-0.426
Female in Urban Area	$\log \hat{u} = 0.465 - 0.420$	89.0	7.0174	-0.420
Female in Rural Area	$\log \hat{u} = 0.907 - 0.547$	84.3	9.0774	-0.547
Nonagricultural Emp.				
Male in Urban Area	$\log \hat{u} = 1.880 - 0.720$	93.0	75.8578	-0.720
Male in Rural Area	$\log \hat{u} = 1.720 - 0.698$	92.8	57.4907	-0.698
Female in Urban Area	$\log \hat{u} = 1.640 - 0.672$	85.2	47.4514	-0.672
Female in Rural Area	$\log \hat{u} = 1.390 - 0.628$	81.8	74.5474	-0.628
Total Employment				
Male in Urban Area	$\log \hat{u} = 2.070 - 0.757$	95.1	117.489	-0.757
Male in Rural Area	$\log \hat{u} = 1.660 - 0.687$	94.5	45.7089	-0.687
Female in Urban Area	$\log \hat{u} = 2.310 - 0.810$	94.1	204.173	-0.810
Female in Rural Area	$\log \hat{u} = 1.430 - 0.639$	95.2		-0.639

The models resulted are evaluated by doing diagnostic analyses to the residual plots. Residual plots for the agricultural employment group are not randomly scattered about the horizontal line  $e=0$  (see Figure 4). Some points are too far away from the line  $e=0$ , since these points have relatively large residuals. The assumption of normality and constant residual variance are also not completely met for each category.

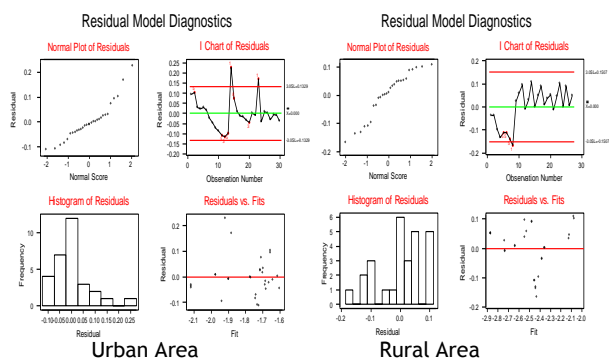
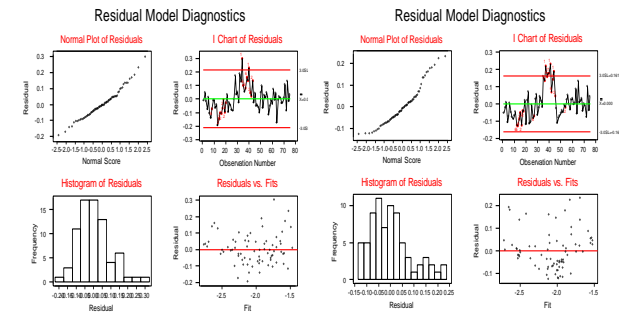


Figure 4. Residual Analyses of the GVF Models for Working Male Category in Agricultural Employment Group

Figure 5 and 6 illustrate that the normal plots of residuals for nonagricultural and total employment group are nearly a straight line. Histograms of residuals are also likely to form a simetric curve, except for the working female categories within the nonagricultural employment

group. Therefore the normality assumption for the residuals is generally met.

often among all households in the sample cluster as well.



Urban Area Rural Area  
 Figure 5. Residual Analyses of the GVF Models for Working Male Category in Nonagricultural Employment Group

A few relatively large residuals in the residual plots for the working female category in urban and rural areas within the nonagricultural employment group may be indicative of outliers-case for which the model is somehow inappropriate. With the exception of a few characteristics with relatively large values of residuals, the GVF models for the nonagricultural employment group (especially for working male category) and the total employment group, have the best empirical performance.

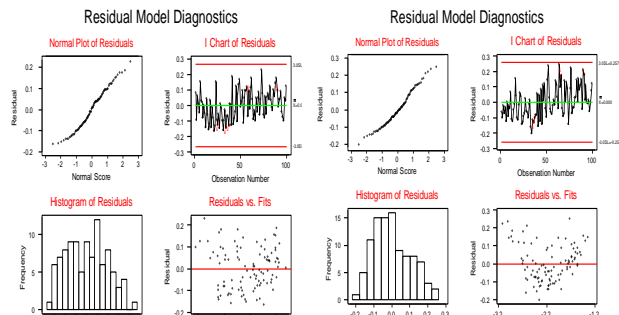
### Illustration for the Application of GVs

By means of the estimated parameters  $c$  and  $d$  in Table 1, the standard error of estimates for working population aged 20-49 years in Java island for which the relative variances are not directly computed, can be calculated. The estimated parameters  $c$  and  $d$  are used according to the interest group of characteristics and demographic categories. For example, the relative variance estimate for the characteristic of working male in urban area as an employee at the industry, is calculated using the values of  $c$  and  $d$  for the category of working male in urban area within the nonagricultural employment group. In Table 1, the parameters are 75.8578 and -0.720, respectively.

Table 2 presents the illustration for calculation of the standard error estimates for each group of characteristics for category of working male in urban area.

Table 2. Illustration for the Application of GVFs for Working Male Category in Urban Area

Group of Employment Sector and Characteristic	Estimate of Total $\hat{Y}$	Model $\hat{u} = c\hat{Y}^d$		Std. Error Estimate $\frac{c}{\hat{Y}^{1-d}}$	95% Confidence Interval	
		C	d		Lower	Upper
(1)	(2)	(3)	(4)	(5)	(6)	(7)
Agricultural Agri. workers- ed. level primary school	592,373	1.8450	0.388	61,053.42	472,708	712,038
Nonagricultural Industry workers- regular employee	1,428,635	75.8578	0.720	75,709.39	1,280,245	1,577,025
Total Ed. level University - aged 35-39	123,139	117.4898	0.757	15,800.87	92,169	154,109



Urban Area Rural Area  
 Figure 6. Residual Analyses of the GVF Models for Working Male Category in Total Employment Group

Based on those empirical results, identification for specific types of characteristics for which GVF will give poor results or inappropriate can be figured out. Those characteristics are the ones that appear to follow a different mathematical model than the rest of the characteristics in the same group. These characteristics have relatively dissimilar deff values to majority of characteristics in the same group. In this case, these characteristics are the common characteristics which usually come out among all persons in the sample household and

### Comparing the Relative Variance Estimates

Figure 7 show a more detailed comparison of GVF and direct computation for working male category in each employment group, respectively. The figure is also a plot of the ratios of relative variance estimates of GVF model to relative variance estimates of direct computation versus the estimate. The findings illustrate that the ratios between the two relative variance estimates range about 0.5 - 1.6, since there are some characteristics having large residuals in the model. It is believed that the use of weighted least squares procedure which prevents characteristics with large relative variances from unduly influencing the estimates of the model parameters, or robust regression method which is less sensitive and more robust to the presence of outliers, will reduce this ratio interval and can improve the model parameter estimation.

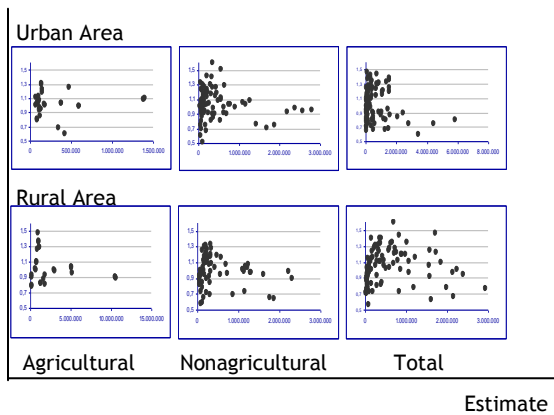


Figure 7. Plot of the Ratios of Relative Variance Estimates of GVF to Relative Variance Estimates of Direct Computation Versus Estimates for Working Male Category in Each Employment Group

### CONCLUSION

The empirical results of simulation study reported here show that model  $\hat{u} = c\hat{Y}^d$  expressed by logarithmic model  $\log \hat{u} = \log c + d \log (\hat{Y})$ , gives a good approximation to estimate the variances for the nonagricultural employment group, especially for working male category both in urban and rural areas. It is also good for the total employment group differentiated by age group, educational attainment and employment status. On the other hand, the model gives poor results for the agricultural employment group.

Based on the empirical results, the GVF models may not perform particularly well for the common characteristics which have relatively dissimilar deff values to majority of characteristics in the same group, since these characteristics usually come out among all persons in the sample household and often among all households in the sample cluster as well.

The success of the GVF technique depends critically on the grouping of the estimates total ( $\hat{Y}$ ) and amount of characteristics involved as the observations for fitting the model. Furthermore, observations with relatively large residuals will also determine the performance of goodness-of-fit of the model.

Application of GVF technique to obtain an approximate standard error on numerous binomial characteristics in large scale survey should be carried out further using extensive data. The better performance of GVF model may also be accomplished by utilizing, for examples, weighted least squares procedure which prevents characteristics with large relative variances from

unduly influencing the estimates of the model parameters, or robust regression method which is less sensitive and more robust to the presence of outliers. Additionally, identification for specific types of characteristics for which GVF will give poor results or may be inappropriate should also be reported.

### REFERENCES

Aunuddin. (1989). *Data Analysis*. Bogor Agricultural University Press, Indonesia.

BPS-Statistics Indonesia. (1998a). *Guidebook for the Supervisor in the 1998 Indonesian Labor Force Survei*. Subdirectorate of Manpower Statistik, BPS-Statistic Indonesia.

BPS-Statistics Indonesia. (1998b). *Labor Force Situation in Indonesia, August 1998*. Subdirectorate of Manpower Statistik, BPS-Statistic Indonesia.

Kish, Leslie. (1995). *Survey Sampling*, New York: John Wiley & Sons.

Minitab Inc. (1996). *Minitab Reference Manual*, Release 11 for Windows. USA.

Smith, P.J. (1998). *Into Statistics : A Guide to Understanding Statistical Concepts in Engineering and the Sciences*, New York: Springer-Verlag.

U.S. Bureau of the Census. (1995). *Cenvar : Variance Calculation System*, IMPS Version 3.1 User's Guide. Washington, D.C.

U.S. Bureau of Labor Statistics. (1996a). *The Current Population Survey: Design and Methodology*, Technical Paper 63, Washington, DC: Government Printing Office.

U.S. Bureau of Labor Statistics. (1996b). *Labor Force Statistics from the Current Population Survey*. Technical Notes to Household Survey Data Published in Employment and Earnings, Washington, DC: Government Printing Office.

Valliant, R. (1987). "Generalized Variance Functions in Stratified Two-Stage Sampling," *Journal of the American Statistical Association*, 82, 499-508.

Weisberg, S. (1985). *Applied Linear Regression*, New York: John Wiley & Sons.

Wolter, K.M. (1985). *Introduction to Variance Estimation*, New York: Springer-Verlag.

