

Probit analysis is applied in a situation where analysis of covariance (ANCOVA) would customarily be used. The dichotomous dependent variables arise from dichotomizations of skewed continuous variables recorded as the proportion of time certain activities are observed. The probit approach avoids the biases of ordinary ANCOVA that arise due to skewness (limited variation). To illustrate this, data from 225 experiments and 214 control subjects in a drug treatment program was analyzed. It was found that the probit approach was able to reveal more substantial treatment effects than the ordinary ANCOVA.

CATEGORIZING SKEWED, LIMITED DEPENDENT VARIABLES

Using Multivariate Probit Regression to Evaluate the California Civil Addict Program

BENGT MUTHÉN
GEORGE SPECKART
University of California, Los Angeles

1. INTRODUCTION

This article is concerned with the problem of multivariate regression analysis where the dependent variables are markedly skewed. In the data to be analyzed, this occurs in the context of variables recorded as the proportion of time a certain activity is observed. Here, the lower limit of zero is observed for a substantial part of the sample, while the upper limit is seldom reached. This situation may be viewed in the context of censored, limited dependent variables, for which it is well known that

AUTHORS' NOTE: *The research of the first author was supported by Grant 81-IJ-CX-0015 from the National Institute of Justice. The research of the second author was supported by National Institute of Drug Abuse Grant DA2577. The authors would like to express their appreciation to M. Douglas Anglin for facilitating access to the data.*

EVALUATION REVIEW, Vol. 7 No. 2, April 1983 257-269
© 1983 Sage Publications, Inc.
0193-841X/83/020257-13\$1.55

ordinary regression analysis will give inconsistent estimates; for example, see Goldberger (1981), Muthén and Jöreskog (1983), and references therein. Ordinary assumptions on the residuals are not realistic in this case. For instance, the conditional expectation of the residual given the exogenous variables cannot be zero, since negative values of the dependent variable cannot be observed.

We will explore the following solution, which uses a probit analysis approach (e.g., see Finney, 1971). We will utilize the concept of a latent response variable underlying each observed dependent variable. The latent variable may be thought of as a response propensity, which when passing thresholds gives rise to different observed outcomes on a categorical variable. Ordinary regression assumptions are assumed to hold for this latent variable as related to the observed exogenous variables. This approach is particularly advantageous when it is believed that little information is lost by categorizing the dependent variable into a small number of categories. We will consider a criminological application where, in the evaluation of a treatment program, strongly skewed response variables such as crime time and incarceration time are observed. In this article, these dependent variables will be dichotomized so that only categories corresponding to "occurrence/nonoccurrence of deviant behavior" will be considered for analysis. For these dichotomous variables, a regression model corresponding to analysis for covariance (ANCOVA) for continuous dependent variables will be formulated. Our methodology draws from Muthén (1979, 1981b), which considered structural equation modeling with a mixture of categorical and continuous latent variable indicators. In our special application, the preferred model to be considered will in fact reduce to the multivariate probit model; see Ashford and Sowden (1970).

In Section 2, the data set is described and the evaluation problem stated. In Section 3, the preferred model is presented. In Section 4, the analyses are carried out, first ignoring the problem using ordinary regression methodology, and then using the preferred model.

2. THE DATA

The data to be utilized in the present analyses were collected by McGlothlin et al. (1977) in order to evaluate the efficacy of treatment of opiate addiction by the California Civil Addict Program (CAP). The

study pioneered by these authors makes a valuable contribution to the knowledge of the benefits of the sometimes controversial approach of civil commitment as a means of treatment and rehabilitation of the drug-abusing criminal offender. The CAP is noteworthy because it represents the largest and most experienced state program of its kind. The appropriateness of utilizing this data for various modeling approaches is underscored by the fact that the only available statistical analyses of this data have simply compared the means of various groups for significant differences.

The first group to be considered consists of a sample admitted to the CAP in 1964 that continued in the normal seven-year commitment. This experimental group is compared to a control group that was discharged from treatment early by writ of habeas corpus. The fortuitous availability of such a control group provides a base level of performance against which to compare the impact of the program on the treatment group. The control group also provides the opportunity to control for otherwise potentially confounding factors that are normally sources of error in most followup studies of this type, such as changes in behavior resulting from increasing age as opposed to treatment effects, biases introduced by errors in recall, or more systematic biases associated with the self-report of sensitive activities such as criminal behavior.

Control subjects were prematched with experimental subjects on the basis of 15 demographic, drug use, and legal status parameters. These included county of commitment, ethnicity, age at first arrest, age at first narcotic use, age at first commitment, marital status, mental status (intelligence), and last grade completed in school. In both the control group and the experimental group, 86% of those selected and not decreased were interviewed in the follow-up period, which occurred some six years after the end of the seven-year commitment for the experimental group. Data will be analyzed for 225 experimental subjects and 214 control subjects for which complete data is available.

Interviewees were paid \$25 for participating; most of the data collected during the interviews assessed changes in legal status, drug use, employment, criminal behavior, and associated monetary estimates of licit and illicit income and expenditures. Further details of the data collection procedures are given in McGlothlin et al. (1977).

For the purposes of the present analyses, the data may be characterized as referring to two separate periods: (1) the precommitment period, defined as that period of time between first narcotic use and admission to the CAP; and (2) the postdischarge period, occurring between the

discharge from the program and the time of the interview. The chief variables of interest for the present analyses are those measuring daily narcotic use, employment, crime, and incarceration. These variables are assigned the labels REGUSE, EMPLOY, CRIME, and INCARC in the tables and discussion that follow. They are additionally differentiated by a "1" or "2" following these labels, designating their reference to either the first or second time periods mentioned above. Furthermore, a "D" affixed at the end of the label signifies a dichotomized or binary variable.

The variables chosen were selected on the basis of their direct relevance to the social costs of drug use. An additional variable was added to the analysis that refers to behavior taking place *during* the commitment period, that is, between time periods 1 and 2. This variable, labeled ABSCOND, indicates the extent to which CAP clients absconded from the supervision of their case officers, and therefore furnishes an estimate of the degree of resistance to the program. Its inclusion was deemed warranted by its predictive utility for postcommitment behavior.

All variables (except those with the D suffix) represent the proportion of time within the associated time period during which the subject was engaged in the designated behavior. Table 1 furnishes frequency counts of the variables from the time 2 or postdischarge period. The categories for which frequencies are reported are numbered 0 to 11, where "0" represents a proportion of 0.0 and "11" designates a proportion of 1.0. The 10 intervening classes each represent a proportion interval of approximately .10 (one-tenth). Accompanying the frequency table is a section listing the mean, standard deviation, skewness, and kurtosis of each variable as computed upon its recoded values. In Table 2 of Section 4 is given the correlation matrix for all variables.

With the exception of the EMPLOY2 variable, the time 2 variables are heavily skewed and may be viewed as censored from below (a large percentage of cases is observed at the lowest value). EMPLOY2 may be viewed as censored both from below and above (bimodal shape). In the preferred model to be presented below, we will dichotomize these variables. To produce a dichotomized version, all variables except EMPLOY2 were recoded to a 0 value if the proportion was 0.0 and to a 1 value otherwise. For EMPLOY2 a 0 value was assigned in the dichotomization if the corresponding proportion was less than or equal to .4999, and a 1 was assigned to values of .50 or higher.

In the control group, those in the 0 category for REGUSE2 (i.e., spent no time using narcotics daily after discharge) consisted of 39%,

TABLE 1
Descriptive Statistics for Dependent Variables (N = 439)

	<i>Frequency, absolute percentage, cumulative percentage for each category code</i>											
	0.	1.	2.	3.	4.	5.	6.	7.	8.	9.	10.	11.
REGUSE2	203	28	36	30	21	22	19	22	13	15	7	23
	46	6	8	7	5	5	4	5	3	3	2	5
	46	53	61	68	72	77	82	87	90	93	95	100
EMPLOY2	57	16	24	24	26	26	31	39	42	41	42	71
	13	4	5	5	6	6	7	9	10	9	10	16
	13	17	22	28	33	39	46	55	65	74	84	100
CRIME2	246	39	35	18	17	15	16	12	10	7	5	19
	56	9	8	4	4	3	4	3	2	2	1	4
	56	65	73	77	81	84	88	91	93	95	96	100
INCARC2	141	40	42	45	35	31	22	27	22	16	14	4
	32	9	10	10	8	7	5	6	5	4	3	1
	32	41	51	61	69	76	81	87	92	96	99	100
	<i>Mean</i>	<i>Standard Deviation</i>					<i>Skewness</i>			<i>Kurtosis</i>		
REGUSE2	2.765	3.466					1.053			-0.147		
EMPLOY2	6.230	3.757					-0.344			-1.197		
CRIME2	2.023	3.156					1.587			1.414		
INCARC2	3.141	3.140					0.716			-0.644		

while in the experimental group the proportion was 53%. For EMPLOY2, the controls in the 0 category (i.e., worked less than half of the time after discharge) were 44%, while experimentals were 35%. For CRIME2, the corresponding proportions for controls and experimentals were 47% and 64% (i.e., those who reported no criminal behavior for time period 2). Finally, the proportions for controls and experimentals regarding the proportion not incarcerated were 26% and 38%, respectively.

Since individuals were not randomly assigned to treatment, the above differences in the outcome variables between experimental and controls cannot safely be taken as treatment effects. Controlling for possible preexisting differences related to these variables seems warranted since they were not directly used in the previously discussed matching procedure. The variable ABSCOND may also indicate uncontrolled pre-treatment differences between experimentals and controls. This leads to a standard ANCOVA where REGUSE1, EMPLOY1, CRIME1, INCARC1, and ABSCOND are used as covariates, and a dummy variable,

TREATMENT, is defined with values 1 and 0 for experimentals and controls respectively. Such an analysis, however, ignores the skewness problem.

3. THE PREFERRED MODEL

In this section, we will specify a multivariate model for the dichotomized variables, which avoids the skewness problem. This model utilizes latent response variables, one for each observed dichotomous variable. For these latent variables will be specified a model analogous to multivariate ANCOVA for continuous dependent variables.

Let \underline{y}_D (4×1) and \underline{x} (6×1) be random vectors of the observed variables, $\underline{y}_D' = (\text{REGUSE2D EMPLOY2D CRIME2D INCARC2D})$, $\underline{x}' = (\text{TREATMENT REGUSE1 EMPLOY1 CRIME1 INCARC1 ABSCOND})$. Let \underline{y}^* (4×1) be a random vector of continuous latent response variables, corresponding to the observed dichotomous variables of \underline{y}_D . In the formulation of Muthén (1981b), each $(\underline{y}^*)_i$ has a dichotomous indicator $(\underline{y}_D)_i$, defined as

$$(\underline{y}_D)_i = \begin{cases} 1, & \text{if } (\underline{y}^*)_i > 0 \\ 0, & \text{otherwise} \end{cases} \quad [1]$$

where i runs over the four dependent variables. It is assumed that \underline{y}^* is related to \underline{x} via the linear structural equation system

$$\underline{y}^* = \underline{\alpha} + \underline{\Gamma} \underline{x} + \underline{\zeta}, \quad [2]$$

where $\underline{\alpha}$ (4×1) is a vector of intercept parameters, $\underline{\Gamma}$ (4×6) is a parameter matrix of regression coefficients, and $\underline{\zeta}$ (4×1) is a random vector of residuals. It is further assumed that $\underline{\zeta}$ has a multivariate normal distribution, is uncorrelated with \underline{x} , has zero expectation, and covariance matrix $\underline{\Psi}$ (4×4). Since the indicator \underline{y}_D of each \underline{y}^* is dichotomous, the scale of \underline{y}^* is indeterminate, and we may standardize the diagonal elements of $\underline{\Psi}$ to unity, expressing off-diagonal elements as residual correlations.

For each of the four dichotomous dependent variables taken separately, this model reduces to a univariate probit model; see Finney (1971). The vector \underline{y}_D follows the multivariate probit model proposed by Ash-

ford and Sowden (1970); these relationships are further described in Muthén (1981a). In the multivariate probit model, both the probit regression coefficients and the residual correlations of $\underline{\Psi}$ are estimated. The multivariate model is a special case of a very general structural equation model put forward in Muthén (1981b). The general model allows for multiple indicators of the latent variables, where a mixture of dichotomous, ordered polytomous, and continuous indicators may be considered. There may also be dependencies among the latent dependent variables. For Case B of Muthén's model no structure is imposed on the exogenous variables. This is the situation of equation 2. With multiple indicators and/or interdependencies among the \underline{y}^* variables, the parameters of the arrays $\underline{\alpha}$, $\underline{\Gamma}$, and $\underline{\Psi}$ of equation 2 obey certain restrictions and are expressed in terms of a smaller set of parameters. For details, see Muthén (1979, 1981a, b).

Ashford and Sowden (1970) considered maximum-likelihood estimation for the multivariate probit model. With more than two dependent variables this method is, however, impractical, since it involves the computation of the multivariate normal distribution function that does not exist in a closed form. In Muthén (1981b) is instead proposed a consistent three-stage estimator. In the first stage each of the dependent variables is regressed on all the exogenous variables by univariate maximum-likelihood probit analysis. In the second stage, all pairs of dependent variables are regressed on all of the exogenous variables by bivariate maximum-likelihood probit analysis. Then the regression coefficients (elements of $\underline{\Gamma}$) are fixed to the estimated values from the first stage so that only a single parameter, the residual correlation (an element of $\underline{\Psi}$) is estimated for each pair of dependent variables. If $\underline{\Gamma}$ and $\underline{\Psi}$ are restricted as discussed above, a third estimation stage is carried out. This is not needed here.

In our application, equation 2 corresponds to multivariate ANCOVA where the first column of $\underline{\Gamma}$ contains treatment effect parameters, and the other columns are the regression coefficients for the covariates. Combining equation 2 with equation 1 defines a new type of ANCOVA for dichotomous dependent variables. As usual, it is implicitly assumed that the coefficients for the covariates are invariant over the treatment and control group, and that group invariance also holds for the residual variances and covariances. With the method of Muthén, such invariance hypotheses can be tested by performing a simultaneous two-group analysis of experimentals and controls. In this article, however, the assessment of invariance is not of primary importance; rather, the comparisons of two techniques given the assumption of invariance is the major issue.

4. THE ANALYSES

We will first analyze the data ignoring the skewness problem. Consider the random vector y' (4×1), $y' = (\text{REGUSE2 EMPLOY2 CRIME2 INCARC2})$. The model is that of equation 2, where we replace the latent vector y^* by the observed vector y . This is then a standard multivariate regression model of the ANCOVA type. This model was estimated by the LISREL computer program (see Jöreskog and Sörbom, 1981), using the maximum likelihood estimator with the correlation matrix. Analyzing the experimental and the control group together gives 439 cases. The Pearson product-moment correlation matrix is given in Table 2.

Table 3 gives the estimated parameters, their standard errors, and their ratios, which will be called test values. Given that the model assumptions are correct, each estimate has an asymptotic normal distribution, so the test values should be compared to the critical values of the z distribution. The results are given for each regression relation in turn; the first three rows for the ordinary regression with the original (skewed) variable, the next three rows for the probit regressions with the corresponding dichotomized variable. At the bottom of the table is given the residual correlation for each relation and each of the types of analysis. For the ordinary regressions with the original variables we note that all significant regression coefficients have the expected signs. All treatment effects are significant, with the largest effect obtained for INCARC2. The R^2 s are all low, and the residual correlations are all moderate to low.

Next we consider the analysis with the preferred model of equations 1 and 2 using the dichotomized variables of y_d . Table 3 gives the regression coefficients from the univariate probit regression, their standard errors, and test values. The estimates are not directly comparable to those of the corresponding ordinary regressions, since different metrics are used. For the regression coefficients, comparisons between the two types of analysis have to be made with respect to test values. The R^2 values for the probit analyses refer to explained variance in the y^* variables, but should be considered as analogous and comparable to the R^2 s of the ordinary regressions. Correlations between the residuals of the four relations are also given and can be compared to those of the ordinary multivariate regression analysis (ANCOVA).

It should be noted that the dependent variable EMPLOY2 was not skewed as were the other three dependent variables and that ordinary

TABLE 2
Correlation Matrix (N = 439)

TREATMENT	REGUSE1	EMPLOY1	CRIME1	INCARC1	ABSCOND	REGUSE2	EMPLOY2	CRIME2	INCARC2
TREATMENT	1.000								
REGUSE1	-0.027								
EMPLOY1	0.081	1.000							
CRIME1	-0.045	-0.263	1.000						
INCARC1	-0.089	-0.220	0.191	1.000					
ABSCOND	0.156	-0.082	0.099	0.076	1.000				
REGUSE2	-0.118	-0.133	0.083	0.099	0.122	1.000			
EMPLOY2	0.116	0.245	-0.052	-0.120	-0.152	-0.480	1.000		
CRIME2	-0.113	-0.114	0.137	0.023	0.142	0.529	-0.418	1.000	
INCARC2	-0.119	-0.147	0.160	0.180	0.251	0.306	-0.328	0.334	1.000

TABLE 3
Model Estimates for Ordinary and Probit Regressions*

Dependent Variable	R ²	Exogenous Variables					
		TREATMENT	REGUSE1	EMPLOY1	CRIME1	INCARC1	ABSCOND
REGUSE2	.075	-.121 .047 -2.57	.165 .048 3.42	-.067 .049 -1.35	-.006 .049 -1.28	.071 .048 1.47	.118 .047 2.51
REGUSE2D	.091	-.399 .124 -3.22	.544 .198 2.75	-.203 .193 -1.05	-.115 .180 -6.39	.198 .360 .055	.547 .233 2.35
EMPLOY2	.100	.117 .046 2.53	-.071 .048 -1.49	.209 .049 4.30	.053 .049 1.09	-.065 .047 -1.37	-.149 .046 -3.21
EMPLOY2D	.108	.287 .127 2.26	-.198 .208 -9.52	.692 .200 3.46	.181 .188 .963	-.321 .365 -8.79	-.637 .251 -2.54
CRIME2	.063	-.132 .047 -2.79	.034 .049 .703	-.074 .049 -1.49	.104 .050 2.09	-.081 .048 -1.67	.149 .047 3.14
CRIME2D	.113	-.471 .125 -3.77	.061 .199 .307	-.292 .192 -1.52	.465 .181 2.57	-.040 .352 -1.14	.524 .246 2.13
INCARC2	.127	-.137 .046 -3.00	.055 .047 1.17	-.057 .048 -1.19	.077 .048 1.61	.124 .047 2.65	.247 .046 5.39
INCARC2D	.227	-.413 .133 -3.11	.546 .219 2.49	-.259 .214 -1.21	-.259 .203 -1.28	1.29 .449 2.87	1.46 .231 6.33

Residual Correlation Matrix							
	REGUSE2	EMPLOY2	CRIME2		REGUSE2D	EMPLOY2D	CRIME2D
EMPLOY2	-.445			EMPLOY2D	-.422		
CRIME2	.513	-.394		CRIME2D	.719	-.460	
INCARC2	.249	-.256	.293	INCARC2D	.560	-.601	.601

*For each dependent variable is given in three different rows, estimates, standard errors, and test values.

regression assumptions may be more closely valid in this case. Dichotomization was still carried out for simplicity and for comparison with skewed cases. The result for the regression is a shrinking of test values when going from EMPLOY2 to EMPLOY2D. This may reflect the loss of information in dichotomization. For the other three dependent variables the changes in test values for the covariates go in both directions. For the treatment effect, which is of primary interest, these three relations all get increased test values for the dichotomized version of the dependent variable. We may note that the relative size of the increases is related to the skewness coefficients reported in Table 1. For CRIME there is a relatively large increase so that with the preferred model, this variable now supersedes INCARC as showing the strongest treatment

effect. All R² values increase when going to the preferred model, reflecting a better specified model. It should be noted that low R² values are to be expected, since the matching variables are not included in the analysis. These variables explain variation in the dependent variables but are not needed as covariates due to the matching procedure. Comparison of the residual correlation matrices shows that dependencies across relations of unexplained variation are in general considerably underestimated in the ordinary regression analysis.

5. CONCLUSION

The methodology proposed here is a special case of a more general formulation presented in Muthén (1981b). This methodology also allows ordered categorical (ordinal) dependent variables with more than two categories, either appearing alone or in mixture with continuous variables. In the preferred model of Section 3, restrictions on the coefficients of Γ and Ψ , such as zeroes, can be readily accommodated by the method. This makes possible the analysis of so-called path analysis models (see Jöreskog and Sörbom, 1981), in this case involving categorical dependent variables.

In this article, we applied a special case of this new methodology to the problem of regression analysis with skewed, limited dependent variables. Such variables are commonly encountered in social science studies where extreme phenomena are observed, such as in criminology and mental health. By categorization (in the present case dichotomization) of such variables, a model can be specified with more realistic statistical assumptions.

We illustrated this approach with an evaluation problem regarding a treatment program with crime time, incarceration time, and narcotics use among the outcome variables. It was found that the preferred methodology of Section 3 resulted in larger estimated treatment effects than the use of ordinary regression analysis. Although treatment effects were all significant as estimated by either method, the larger effects yielded by the preferred method can be important with smaller samples or smaller true effects; a true effect might otherwise go unnoticed.

Other analyses approaches for this kind of data should be noted. In one alternative, we may dichotomize the dependent variables as above and perform ordinary ANCOVA. For this data set, analyses suggest that smaller differences in treatment effects occur when comparing this approach to the one presented in Section 3. That is, we obtain a larger

distortion when using ordinary methods on skewed continuous (12 category) dependent variables, than with dichotomized variables. Although ordinary methods on dichotomized variables give a smaller distortion, the Section 3 methodology is still preferable. As with ordinary probit (logit) regression compared to ordinary regression with a dichotomous dependent variable, this is particularly true with more extreme dichotomizations (proportions less than .2, say).

Another alternative is to analyze the original continuous (12 category) skewed variables by Tobit modeling (see Tobin, 1958; Amemiya, 1982). This modeling explicitly takes into account the limited variation of the variable, viewing it as censored from below.

Still another approach is taken in Muthén and Speckart (1982). This further develops the modeling of Section 3 to using a single latent dependent variable ("deviant behavior"), being measured by the four dichotomized variables using a dichotomous factor analysis measurement model (see Muthén, 1979). Treatment effects are then studied in the single regression of the latent variable on the exogenous variables. When several outcomes measure more or less the same thing, such modeling can be attempted, and if well fitting, more power will result in determining treatment effects.

REFERENCES

- AMEMIYA, T. (1982) Tobit Models: A Survey. CJRS 4. Palo Alto, CA: Rhodes Associates.
- ASHFORD, J. R. and R. R. SOWDEN (1970) "Multivariate probit analysis" *Biometrics* 26: 535-546.
- FINNEY, D. J. (1971) *Probit Analysis*. Cambridge: Cambridge Univ. Press.
- GOLDBERGER, A. S. (1981) "Linear regression after selection" *J. of Econometrics* 15: 357-366.
- JÖRESKOG, K. G. and D. SÖRBOM (1981) LISREL V—Analysis of Linear Structural Relationships by Maximum-Likelihood and Least Squares Methods. Research Report 81-8, Department of Statistics, University of Uppsala.
- McGLOTHLIN, W. H., M. D. ANGLIN and B. D. WILSON (1977) *An Evaluation of the California Civil Addict Program*. NIDA Services Research Monograph Series, Rockville, Maryland.
- MUTHÉN, B. (1981a) "Some categorical response models with continuous latent variables," in the conference volume, K. G. Joreskog and H. Wold (eds.) *Systems Under Indirect Observation: Causality, Structure, Prediction*. Amsterdam: North-Holland.
- (1981b) *A General Structural Equation Model with Ordered Categorical and Continuous Latent Variable Indicators*. Department of Psychology, University of California, Los Angeles.

- (1979) "A structural probit model with latent variables." *J. of the Amer. Statistical Ass.* 74: 807-811.
- and K. G. JÖRESKOG, (1983) "Selectivity problems in quasi-experimental studies" *Evaluation Review*.
- MUTHÉN, B. & G. SPECKART (1982) *Latent Variable Probit ANCOVA: Treatment Effects in the California Civil Addict Program*.
- TOBIN, J. (1958) "Estimation of relationships for limited dependent variables." *Econometrica* 26: 24-36.

Bengt Muthén is Assistant Professor of Education at the Graduate School of Education, UCLA. Dr. Muthén was formerly a researcher in the Joreskog group of the Statistics Department, University of Uppsala, where he earned his Ph.D in Statistics. Dr. Muthén's research involves latent variable structural equation modeling with categorical data, including the factor analysis of dichotomous variables.

George Speckart is a psychology graduate student at UCLA. He is currently working on a doctoral dissertation utilizing structural equation models to investigate the etiological relationship between property crime and narcotics use. During the last four years he has worked at UCLA with Dr. Williams McGlothlin collecting follow-up data from methadone patients.