

ADDIS ABABA UNIVERSITY
SCHOOL OF GRADUATE STUDIES

TITLE OF RESEARCH DETERMINANTS OF FERTILIZER USE IN THE NORTH-WEST

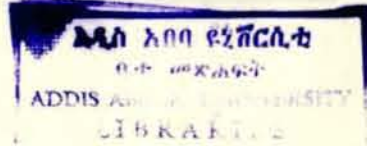
DETERMINANTS OF FERTILIZER USE IN THE NORTH-WEST

ETHIOPIA

AN APPLICATION OF QUALITATIVE RESPONSE MODEL

SCIENCE

By



APPROVED BY

Desta Woldemariam

Internal Examiner

A THESIS SUBMITTED IN PARTIAL FULFILLMENT

FOR THE DEGREE MASTER OF SCIENCE IN STATISTICS

IN THE ADDIS ABABA UNIVERSITY

Addis Ababa

June 1990

ADDIS ABABA UNIVERSITY
SCHOOL OF GRADUATE STUDIES

TITLE OF RESEARCH DETERMINANTS OF FERTILIZER
USE IN NORTHERN ETHIOPIA
AN APPLICATION OF QUALITATIVE RESPONSE MODELS

NAME OF CANDIDATE DESTA WOLDEMARIAM

DEPARTMENT STATISTICS

FACULTY SCIENCE

APPROVED BY

Advisor & Chairman

External Examiner

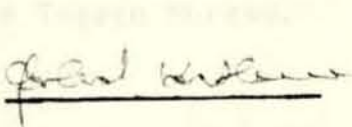
Internal Examiner

Name Asmerom Kidanu

GEORHARD KOCKLÄNDER

IGOR LITVIN

Sign. 





Date JUNE 14, 1990

ACKNOWLEDGEMENTS

I am highly indebted to Dr. Asmerom Kidane, my Advisor; for his unreserved advice, support in data collection and his personal computer facility. Had it not been for his help, this research work could have been incomplete. Financial grant from Rockefeller Foundation is highly appreciated. I am also grateful to W/t Zenebech W. Tsadik for spending her valuable time to type the manuscript, and particularly I acknowledge Ato Tegegn Nuresu.

ABSTRACT

In this paper the method of Qualitative Response Models (QRM) is applied to investigate the determinants of fertilizer usage among Ethiopian farmers in the North-west. The theoretical basis of the decision is based on the concept of stochastic utility model. A survey data from rural Ethiopia was applied to test the model.

The application of linear probability, probit and logit models seem to be satisfactory in explaining variability in fertilizer usage. The study showed that most of the explanatory variables (farm land size, the household's head education level, number of farm implements, number of cattle, number of plots, number of males aged 10-14, number of females aged 10-14, number of males aged 15-54, number of females aged 15-54 and non-farm income) affect positively the probability that a farmer will use fertilizers. Probabilities using minimum, mean and maximum values of the explanatory variables are also generated.

Finally, five selected non-demographic variables, namely, the level of education of the households head, number of farm implements, farm land size, number of cattle, number of implements and non-farm income were increased by various rates from their mean values, while demographic variables remained fixed at their mean values. This helped

CONTENTS

us to generate various probabilities. The result was an increase in the farmer's probability of using fertilizers. This suggests that measures have to be taken so as to make the above incentives practical. 1

 1.2 Steps of the study 2

CHAPTER 2 BASIC THEORETICAL CONCEPTS 4

 2.1 Qualitative response models 4

 2.2 Estimation and hypotheses testing 10

 2.2.1 Maximum likelihood (ML) estimator and minimum chi-square (min χ^2) estimator for QDM 12

 2.2.2 Iteration 17

 2.2.3 Tests of hypotheses 14

 2.3 Choice of models 16

 2.4 Multi-rasour models 19

 2.5 Unordered independent logit model 20

CHAPTER 3 SOME APPLICATIONS OF QUALITATIVE RESPONSE MODELS 23

CHAPTER 4 DATA 37

CHAPTER 5 APPLICATION 33

 5.1 Utility maximization theory 33

 5.2 Variables 36

 5.3 Estimation methods and results .. 37

 5.4 DISCUSSION 41

CONTENTS

	<u>PAGE</u>
CHAPTER 1 INTRODUCTION.....	1
1.1 General	1
1.2 Steps of the study	2
CHAPTER 2 BASIC THEORETICAL CONCEPTS	4
2.1 Qualitative response models	4
2.2 Estimation and hypotheses testing	10
2.2.1 Maximum likelihood (ML) estimator and minimum chi-square (min χ^2) estimator for QRM	12
2.2.2 Iteration	12
2.2.3 Tests of hypotheses	14
2.3 Choice of models	16
2.4 Multi-response models	19
2.5 Unordered independent logit model	20
CHAPTER 3 SOME APPLICATIONS OF QUALITATIVE RESPONSE MODELS	23
CHAPTER 4 DATA	32
CHAPTER 5 APPLICATION	33
5.1 Utility maximization theory	33
5.2 Variables	36
5.3 Estimation methods and results ..	37
5.4 Discussion	41

	<u>PAGE</u>
5.5 Probability generation	42
5.6 Comparison of models	47
5.7 Tests of hypotheses	47
CONCLUSIONS	51
APPENDIX: A. Table of logit model estimates	54
B. Table of probit model estimates ...	55
C. Table of linear probability model estimates	56
REFERENCES.....	57

This study will deal with the last point, that is, the case of fertilizers used present farms. More specifically we will try to identify the important determinants in the demand for fertilizers. In other words we will try to identify variables that motivate farmers to use more or less of no fertilizers.

If we are able to identify farmers that choose farmers to use fertilizers, then the way to go is to provide the

I. INTRODUCTION

11. General

It is widely known that over 90 percent of the people in Ethiopia are engaged in agriculture and this sector is the mainstay of Ethiopian economy. Besides, over 80% of foreign exchange earnings come from this sector. However, agricultural productivity in Ethiopia is below expectation. In developed countries less than 10% of the population are engaged in agriculture and produce not only for domestic consumption but also for export.

This is not the case for the less developed countries (LDC). Ethiopia is a case in point. The country is deficient in food production; drought and famine seem to be the order of the day. Cognizant of this fact there is an attempt to increase agricultural productivity through various forms of assistance such as the provision of better quality seeds and fertilizers.

This study will deal with the last point, that is, the case of fertilizers among peasant farms. More specifically we will try to identify the important determinants in the demand for fertilizer. In other words we will try to identify variables that motivate farmers to use more or less or no fertilizer.

If one is able to identify factors that induce farmers to use fertilizer, then one may be able to suggest to policy

makers on the efficient distribution of fertilizers. This will result in efficient utilization of scarce foreign exchange and also may help increase the productivity of peasant agriculture. The benefit would be self sufficiency in food production as well as possibilities of exporting surplus food.

The method of analysis will be the application of Qualitative Response Models (ORM). Among many, emphasis will be put on three common methods of analysing qualitative or categorical data, namely the linear probability, probit and logit models.

1.2 Steps of the Study

The study will have four parts: First ^{we} will consider the properties of the various qualitative response models and their relevance in measuring the type of response we are trying to study. This will help us generate probabilities of using or not using fertilizer under various socioeconomic and environmental conditions.

Secondly, we will review various application of qualitative response models to the data.

Thirdly the economic rationale involved in farmers decisions making process will be discussed by building appropriate mathematical approach of stochastic utility models. This will help us to model a peasants' behaviour under different

socioeconomic, environmental conditions and other constraints.

Fourth; we will identify appropriate explanatory variables and provide reasons for the inclusion of some. We will then apply the three qualitative response models and generate appropriate probabilities by varying the explanatory variables.

The data for this research is obtained from a socioeconomic and demographic multipurpose sample survey in the north/west region of Ethiopia collected during 1989, financial grant by Rockefeller Foundation for the research work on "Economic and Demographic Household Behavior in Rural Ethiopia."

A stratified multistage cluster sampling was applied to select about 800 households from the region and finally 661 observations are considered here.

The data is processed using computer and applying a statistical software STATA.

Let y be a dichotomous random variable which takes the value 1 if the event occurs or 0 if it does not, (though any other pair of real numbers could be used, the choice of 0 and 1 is especially convenient) (2).

We assume that the probability of an event depends on a vector of independent variables X and a vector of unknown parameters β . Using subscript i to denote the i^{th} individual, we can write a univariate distribution model as

$$P_i = \frac{\exp(\beta_0 + \beta_1 X_i)}{1 + \exp(\beta_0 + \beta_1 X_i)} \quad (2.1)$$

II. BASIC THEORETICAL CONCEPTS

2.1 Qualitative Response Models

Distinction has to be taken between qualitative and quantitative data. In the former case observations are recorded by proxies while in the latter case observations are recorded in figures.

One of the most important development in econometrics in the past has occurred in the area of qualitative response models, abbreviated as QRM. Also known as quantal, categorical or discrete models. Here the dependent variable is qualitative or categorical.

Suppose we want to consider the occurrence and non-occurrence of an event such as "a farmer uses fertilizer or not" in our case, it is mathematically convenient to define a dichotomous random variable y which takes the value 1 if the event occurs or 0 if it does not, (though any other pair of real numbers could be used, the choice of 0 and 1 is especially convenient) [2].

We assume that the probability of an event depends on a vector of independent variables X and a vector of unknown parameter θ . Using subscript i to denote the i^{th} individual, we can write a univariate dichotomous model generally as

$$p_i = p(y_i=1) = G(X_i;\theta) \quad (2.1)$$

$$i = 1, \dots, n$$

where y_i is defined as:

$$y_i = \begin{cases} 1 & \text{if the event occurs} \\ 0 & \text{Otherwise} \end{cases}$$

Here we assume that y_i 's are independent. Equation (2.1) states that, in our case the probability of fertilizer usage depends on the farmers socio-economic, demographic and characteristics vector X_i .

There are many OR models, but the most frequently used three common models are linear probability, probit and logit models [2].

(2.2) Linear probability (LP) model:

$$F(w) = w$$

probit model:

$$F(w) = \Phi(w) = \int_{-\infty}^w \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$$

Logit model:

$$F(w) = L(w) = \frac{e^w}{1+e^w}$$

For $v = \chi_i' \beta$, the linear probability model has an obvious defect in that $\chi_i' \beta$ is not constrained to lie between 0 and 1 as a probability should. Though this defect can be corrected by defining

$$(2.3) \quad F = \begin{cases} 1 & \text{if } \chi_i' \beta > 1 \\ \chi_i' \beta & \text{if } 0 \leq \chi_i' \beta \leq 1 \\ 0 & \text{if } \chi_i' \beta < 0 \end{cases}$$

Since $E(y_i) = p(y_i = 1)$, we have in Ln model.

$$(2.4) \quad y_i = x_i' \beta + U_i$$

where $u_i = y_i - x_i' \beta = y_i - E(y_i)$. This is a heteroscedastic regression model since $V(U_i) = V(y_i) = x_i' \beta (1 - x_i' \beta)$, using the variance formula for a binomial variable, provided $0 < \beta < 1$. The Least Square (LS) method yields consistent and unbiased estimates of β and the weighted least squares (WLS) method yields consistent and asymptotically more efficient estimates of β . However, neither LS nor WLS method avoid the inherent weakness of the model, heteroscedasticity. The WLS procedure fails if the condition $0 < x_i' \beta < 1$ is not met. It is better to use LS rather than WLS. However, we should be aware the fact that the standard deviations of the LS estimates given by the standard LS program are biased because of the heteroscedasticity [2].

The probability function used for the probit model is the standard normal distribution function, and the logit model has used the logistic distribution function. Both distribution functions are bounded between 0 and 1, and are symmetric around 0. A random variable from the probit model and logit model has variance 1 and $\pi^2/3$ respectively.

If we consider a transformed logistic distribution such as

$$(2.5) \quad L\lambda(w) = \frac{e^{\lambda w}}{1 + e^{\lambda w}}$$

We can make the model to closely approximate the standard normal distribution by choosing an appropriate value of λ over a wide domain. Amenyia (1981) has shown that $\lambda = 1.6$ i.e. $L_{1.6}$ is a good approximate.

Results from the two models (probit and logit) are similar. Because of the close similarity of the two distribution, it is difficult to distinguish between them statistically unless one has extremely large number of observations (chambers and Cox, 1976). Thus in the univariate dichotomous model, it does not matter much whether one uses a probit model or a logit model, except in cases where data are heavily concentrated in the tails due to the characteristics of the problems studied [2].

As shown in Amenyia (1981), if $\hat{\beta}$ is the estimate of β , then approximately

$$(2.6) \quad 1.6 \hat{\beta}_{\phi} \approx \hat{\beta}_L$$

where $\hat{\beta}_{\phi}$ is estimate of the coefficient in the probit model and $\hat{\beta}_L$ of the logit model. This formula is useful if one is needed a quick way to compare probit and logit models. In general

$$(2.7) \quad \hat{\beta}_{L^p} = 0.4 \hat{\beta}_{\phi}$$

where $\hat{\beta}_{L^p}$ = estimate of β for the L^p model (except the constant term) $\hat{\beta}_{L^p} = 0.4\hat{\beta}_{\phi} + 0.5$ for the constant term, and

$$(2.8) \quad \begin{aligned} \hat{\beta}_{L^p} &= 0.25 \hat{\beta}_L \quad (\text{except for the constant term}) \\ \hat{\beta}_{L^p} &= 0.25 \hat{\beta}_L + 0.5 \quad \text{for the constant term.} \end{aligned}$$

The above gives us a quick and rough approximation. However, to compare probability functions of different models, it is generally better to compare probability directly rather than comparing estimates of coefficients even after an appropriate conversion [2].

An alternative way of comparing different models is to look at the derivative of the probabilities with respect to a particular independent variable. Let x_{ik} be the k^{th} element of x_i and β_k be the k^{th} element of β . Then, the derivatives for the three probability models are given by

$$(2.9) \quad \begin{aligned} \frac{\partial x_i' \beta}{\partial x_{ik}} &= \beta_k \\ \frac{\partial \phi(x_i' \beta)}{\partial x_{ik}} &= \phi(x_i' \beta) \beta_k \\ \frac{\partial L(x_i' \beta)}{\partial x_{ik}} &= \frac{e^{x_i' \beta}}{(1 + e^{x_i' \beta})^2} \cdot \beta_k \end{aligned}$$

$$i = 1, \dots, n$$

where ϕ is the standard normal density function. If, in the above, the right-hand side evaluated at $x_i' \beta = 0$, (2.7) and (2.8) will be approximately obtained.

As it is explained in Amemya (1981), similarity between the probit ML and LP-WLS estimates is noted by Hill (1979) and a similarity between the logit ML and Ln-LS estimates is noted by Wilensky and Rossiter (1978).

Now let us consider how QR models are specified in economic applications. Economists assume that an economic

unit makes rational decisions so as to maximize its utility.

Suppose a consumer has two alternatives to satisfy his need, say modes (mode 0 and 1) and suppose the utility is a function of the mode characteristics z , individuals socio-economic characteristics w , plus an additive error term e . We define U_{i0} and U_{i1} as the i^{th} person indirect utilities associated with mode 0 and mode 1 respectively. Assuming linear function we have

$$(2.10) \quad \begin{aligned} U_{i0} &= \alpha_0 + z_{i0}'\beta + w_i'\gamma_0 + \epsilon_{i0} \\ U_{i1} &= \alpha_1 + z_{i1}'\beta + w_i'\gamma_1 + \epsilon_{i1} \end{aligned}$$

utility function can be presented for the general case as

$$U_{ij} = \mu_{ij} + \epsilon_{ij} \quad \begin{aligned} i &= 1, \dots, n \\ j &= 0, 1 \end{aligned}$$

where μ_{ij} is non-stochastic function of the explanatory variables.

The basic assumption is individual i chooses model 1 if $U_{i1} > U_{i0}$ and model 0 will be chosen if $U_{i0} > U_{i1}$. If we assume ϵ_{i1} and ϵ_{i0} to be continuous random variables and define

$$(2.11) \quad \begin{aligned} p(y_i=1) &= p(U_{i1} > U_{i0}) \\ &= p(\mu_{i1} + \epsilon_{i1} > \mu_{i0} + \epsilon_{i0}) \text{ general case} \\ &= p(\epsilon_{i0} - \epsilon_{i1} < \mu_{i1} - \mu_{i0}) \\ &= F(\mu_{i1} - \mu_{i0}) \end{aligned}$$

where F is the distribution function of $\epsilon_{i0} - \epsilon_{i1}$.

Therefore, what kind of QR models one gets is equivalent to what distribution one assumes to $\epsilon_{i0} - \epsilon_{i1}$. For example a probit (logit) model arises from assuming normal (logistic) distribution for $\epsilon_{i0} - \epsilon_{i1}$.

2.2 Estimation and Hypothesis Testing

Let y be a dichotomous dependent variable and x be a vector of independent variables. We will associate the event that the vector (y, x') takes on a particular vector value with the word cell.

The estimation of QR model is simpler if there are many observations (≥ 30) per cell. The case of many observations per cell doesn't occur (appear) frequently, particularly in economic applications [2].

2.2.1 Maximum Likelihood (ML) Estimator and Minimum Chi-Square (Min χ^2) Estimator for QR Models

The ML estimator can be used either in the case of few observations per cell or many observations per cell, where as the minimum χ^2 can be effectively used only in the case of many observations per cell.

(i) Few observations per cell

ML estimator

The maximum likelihood function is given by

$$(2.12) \quad LF = \prod_{i=1}^n F(x_i' \beta)^{y_i} [1 - F(x_i' \beta)]^{1-y_i}$$

where $y_i = 0$ or 1

and its natural logarithm

$$(2.13) \quad \ell = \sum_{i=1}^n y_i \log F(x_i' \beta) + \sum_{i=1}^n (1-y_i) \log [1 - F(x_i' \beta)]$$

The ML estimator $\hat{\beta}_{ML}$ is defined as the value of β that maximizes either (2.12) or (2.13). Differentiating ℓ with respect to the column vector β yields a column vector of derivatives.

$$(2.14) \quad \frac{\partial \ell}{\partial \beta} = \sum_{i=1}^n \frac{y_i - F(x_i' \beta)}{F(x_i' \beta) [1 - F(x_i' \beta)]} f(x_i' \beta) x_i$$

where f denotes the derivative of F . In case of probit $F = \Phi$, $f = \phi$ and in case of the logit $F = L$, $f = L(1-L)$. $\hat{\beta}_{ML}$ is a solution of the equation.

$$(2.15) \quad \frac{\partial \ell}{\partial \beta} = 0$$

Since (2.15) is non-linear in β , the ML estimator has to be obtained by an iterative method. In order to define an iterative method as well as to derive the asymptotic variance - covariance matrix, we need the second order derivatives of ℓ and their expectation. Differentiating the

column vector $\partial \ell / \partial \beta$ with respect to the row vector β yields a matrix of second order derivatives.

$$(2.16) \quad \frac{\partial^2 \ell}{\partial \beta \partial \beta'} = - \sum_{i=1}^n \left[\frac{y_i}{F^2(x_i' \beta)} - \frac{1-y_i}{[1-F(x_i' \beta)]^2} \right] f^2(x_i' \beta) x_i x_i'$$

$$(2.20) \quad \text{Method of } + \sum_{i=1}^n \left[\frac{y_i - F(x_i' \beta)}{F(x_i' \beta) [1-F(x_i' \beta)]} \right] f'(x_i' \beta) x_i x_i'$$

where f' is the derivative of f . Taking the expectation of (2.16) we get

$$(2.17) \quad E \left\{ \frac{\partial^2 \ell}{\partial \beta \partial \beta'} \right\} = - \sum_{i=1}^n \frac{f^2(x_i' \beta)}{F(x_i' \beta) [1-F(x_i' \beta)]} x_i x_i'$$

(2.21) It is well known that under general conditions the Newton-Raphson maximum likelihood estimator is consistent and asymptotically normal with variance covariance matrix $-(E \partial^2 \ell / \partial \beta \partial \beta')^{-1}$.

(See Dobson, 1983). Therefore

$$(2.18) \quad V(\hat{\beta}_{ML}) = \left[\sum_{i=1}^n \frac{f^2(x_i' \beta)}{F(x_i' \beta) [1-F(x_i' \beta)]} x_i x_i' \right]^{-1}$$

where An estimate of $V(\hat{\beta}_{ML})$ is obtained by evaluating (2.18) at $\hat{\beta}_{ML}$.

2.4 Iteration

The two most commonly used iterative methods for calculating the ML estimator are "Newton-Raphson" method and

the method of "scoring". Given an initial estimate $\hat{\beta}_1$, the second round estimate $\hat{\beta}_2$ in each method is defined as follows:

(2.19) Newton-Raphson method:

$$\hat{\beta}_2 = \hat{\beta}_1 - \left[\frac{\partial^2 \ell}{\partial \beta \partial \beta'} \Big|_{\beta_1} \right]^{-1} \frac{\partial \ell}{\partial \beta} \Big|_{\beta_1}$$

(2.20) Method of scoring:

$$\hat{\beta}_2 = \hat{\beta}_1 - \left\{ E \left[\frac{\partial^2 \ell}{\partial \beta \partial \beta'} \Big|_{\beta_1} \right] \right\}^{-1} \frac{\partial \ell}{\partial \beta} \Big|_{\beta_1}$$

The third-round, fourth-round, etc. are to be evaluated using the above procedure. General formula can be given for each iteration method (Dobson, 1983) [11].

(2.21) Newton-Raphson

$$b^m = b^{m-1} - \left[\frac{\partial^2 \ell}{\partial \beta_j \partial \beta_k} \Big|_{\beta=b^{(m-1)}} \right]^{-1} U^{(m-1)}$$

(2.22) Method of scoring:

$$b^{(m)} = b^{(m-1)} + [I^{(m-1)}]^{-1} U^{(m-1)}$$

where

$$U = [U_1, \dots, U_p]' \quad , \quad I^{(m-1)} = - \left\{ F \frac{\partial^2 \ell}{\partial \beta \partial \beta'} \Big|_{\beta=b^{(m-1)}} \right\}$$

$$U_j = \frac{\partial \ell}{\partial \beta_j} \quad , \quad j=1, \dots, p$$

I = information matrix

$b^m = m^{\text{th}}$ round estimate of β .

The method of scoring can be rewritten as

$$(2.23) \quad \hat{\beta}_2 = \left[\sum_{i=1}^n \frac{\hat{f}_i^2}{\hat{F}_i(1-\hat{F}_i)} x_i x_i' \right]^{-1} \sum_{i=1}^n \frac{\hat{f}_i}{\hat{F}_i(1-\hat{F}_i)} x_i' [y_i - \hat{F}_i + \hat{f}_i' \epsilon_1]$$

where $\hat{F}_i = F(x_i' \hat{\beta}_1)$ and $\hat{f}_i = f(x_i' \hat{\beta}_1)$

and from (2.4) $y_i = F(X_i' \beta) + U_i$ where $E(U_i) = 0$ and $V(U_i) = F(X_i' \beta)(1-F(X_i' \beta))$.

(ii) Many observation per cell

ML methods as well as Min χ^2 method can be applied. We will not consider here the case of many observation per cell.

2.5 Tests of Hypothesis

Tests of a general linear hypothesis of the form

$$(2.24) \quad Q' \beta = C$$

where Q' is $m \times k$ matrix of known constants (k being the number of elements of β) and C is an m vector of known constants. Assuming that $m \leq k$ and m rows of Q' are linearly independent.

For $m=1$, to test the hypothesis

$$H_0: Q' \beta = C$$

the test is based on

$$(2.25) \quad \frac{Q' \hat{\beta} - C}{\sqrt{Q' (\hat{V} \hat{\beta}) Q}} \stackrel{A}{\sim} N(0,1)$$

where $\hat{V} \hat{\beta}$ = consistent estimator of the asymptotic variance $V \hat{\beta}$ and A stands for asymptotically.

$$(2.26) \quad \frac{C'\hat{\beta} - C}{\sqrt{C'(\hat{V}\hat{\beta})Q}} \overset{A}{\sim} t_{n-k}$$

where t_{n-k} denotes student's t distribution with $n-k$ degrees of freedom.

For $m > 1$, Two tests will be considered

- i) Wald's test
- ii) The likelihood ratio test (LRT).

Wald's test can be used in connection with any estimator, where as the LRT must be based on either the ML estimator or any estimator with the same asymptotic distribution. In using these tests we must always assume that the alternative hypothesis $C'\beta \neq C$. Though these tests are valid even when $m=1$.

$$(2.27) \quad \text{Wald} = (C'\hat{\beta} - C)' [C'(\hat{V}\hat{\beta})Q]^{-1} (C'\hat{\beta} - C) \overset{A}{\sim} \chi_m^2$$

The hypothesis is to be rejected if the value of the statistic exceeds a prescribed critical value.

The likelihood ratio is defined by

$$(2.28) \quad \text{LRT} = 2[\ell(\hat{\beta}_{ML}) - \ell(\hat{\beta}_{CML})] \overset{A}{\sim} \chi^2$$

where $\hat{\beta}_{CML}$ denotes the constrained maximum likelihood (CML) estimator.

2.3 Choice of Models

It is noted that the probit and logit models usually give similar results and it is difficult to distinguish them statistically [2]. However, Kimio Morimune (1979) [22] has suggested the following in his concluding remarks of research work on "comparison of normal and logistic models in the bivariate dichotomous analysis."

The original objective was to compare relative performances of three probability models: the logistics, normal and linear models. However, the linear model was found to be far inferior to the other two models. Then the research is confined to the comparison of the normal and logistic models.

Using Cox-test, the experiment shows that the standard error of estimated coefficients of the normal (probit) model tend to be smaller than those of the logistic model. This has also been supported by some numerical comparisons of generalized variance of estimators. However, more research is required to see this aspect of comparison in detail. Some criterion may help us in choosing among competing models. Some of the criterion are R^2 , number of wrong prediction, sum of squared residuals (SSR), squared correlation coefficient, etc.

Except the squared correlation coefficient and the log likelihood function, we have to choose the model for which the value of the criterion is smallest. In the passage, \hat{F}_i

denotes $F(x_i, \hat{\beta})$ where $\hat{\beta}$ is whatever estimator is being used.

(2.28) - number of wrong predictions:

$$(2.29) \quad \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

where
$$\hat{y}_i = \begin{cases} 1 & \text{if } \hat{F}_i > \frac{1}{2} \\ 0 & \text{if } \hat{F}_i < \frac{1}{2} \end{cases}$$

The value gives the number of wrong predictions because $(y_i - \hat{y}_i)^2 = 1$ if and only if $y_i \neq \hat{y}_i$.

A major disadvantage is that if we are dealing with an event which happens with a high probability (eg. a man working) or a low probability (eg. a person immigrating), most models will do well by this criterion [2].

Sum of squared residuals (SSR):

$$(2.30) \quad \sum_{i=1}^n (y_i - \hat{F}_i)^2$$

This criterion does not suffer from the deficiency of the number of wrong predictions. This is a natural criterion, since it corresponds to the sum of squared residuals in the standard regression model, from which R^2 is derived. However, its use in OR models cannot be defined as strong as in the standard regression model because a QR model is essentially a heteroscedastic regression model [1]. Efron (1978) defends (2.30) from a certain axiomatic point of view and suggests an analogue of R^2 defined by

$$(2.31) \quad \text{Bfren's } p^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{F}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

where
$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

- SSR weighted by estimated probabilities:

$$(2.32) \quad \sum_{i=1}^n \frac{(y_i - \hat{F}_i)^2}{\hat{F}_i(1-\hat{F}_i)}$$

For two reasons we prefer this criterion to the un-weighted SSR. First, if the true probabilities were known and used in the denominator instead of the estimated probabilities, the minimization of the above criterion with respect to β yields the estimator of β which is asymptotically more efficient than that obtained by the minimization of the un-weighted SSR. Second, it seems reasonable to attach a higher loss to the error made in predicting a random variable with a smaller variance since such a random variable should be easier to predict than the one with a larger variance. Thus, it seems reasonable to weigh the squared error by a weight which is inversely proportional to the variance.

- Squared correlation coefficient:

$$(2.33) \quad \frac{\left[\sum_{i=1}^n (y_i - \bar{y}) \hat{F}_i \right]^2}{\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{F}_i - \bar{F})^2}$$

This measure is closely related to the un-weighted SSR (230), the same criticism applies to the squared correlation coefficient as to SSR.

- Log likelihood function:

$$(2.34) \quad \ell = \sum_{i=1}^n [(y_i \log \hat{F}_i + (1-y_i) \log (1-\hat{F}_i))]$$

Where \hat{F}_i here specifically denotes $F(x_i, \hat{\beta}_{ML})$. This measure has an obvious **intuitive** appeal. In addition, it is especially suitable for comparison of different numbers of parameters. Suppose we want to choose between the unconstrained model in which the k-component parameter vector β is allowed to vary freely and the constrained model in which β is subject to q linear constraints $C'\beta=C$, (q = number of constraints). One should choose the unconstrained model if and only if $2[\ell(\hat{\beta}_{ML}) - \ell(\hat{\beta}_{OML})]$ is greater than the $\alpha\%$ critical value of χ^2_q .

2.4 Multi-response Models

Assuming that the dependent variable y_i takes m_i+1 values $0, 1, 2, \dots, m_i$, a general multi-response QR model can be written as

$$(2.35) \quad p(y_i = j) = F_{ij}(x^*, \theta)$$

$$i = 1, 2, \dots, n \text{ and}$$

$$j = 0, 1, 2, \dots, m_i$$

where x^* and θ are vectors of independent variables and parameters respectively.

Defining $\sum_{i=1}^n (m_i + 1)$ binary variables

$$(2.36) \quad y_{ij} = \begin{cases} 1 & \text{if } y_i = j \\ 0 & \text{if } y_i \neq j \end{cases}$$

$i = 1, 2, \dots, n$ and

$j = 0, 1, \dots, m_i$

The likelihood function of model (2.35) is given as

$$(2.37) \quad LF = \prod_{i=1}^n \prod_{j=0}^{m_i} p_{ij}^{y_{ij}}$$

General results about the asymptotic distribution and the iterative methods concerning ML estimation discussed previously apply for the present model.

2.5 Unordered Independent Logit Model

By specifying the probability function (2.35) for a certain i for which $m_i = 2$ (the case of larger m_i will be inferred from the following). Writing

$$p_{ij} = p(y_i = j), \quad j = 0, 1 \text{ and } 2,$$

the three probabilities are specified by

$$(2.38) \quad p_{i2} = \frac{e^{X'_{i2}\beta}}{1 + e^{X'_{i1}\beta} + e^{X'_{i2}\beta}}$$

$$(2.39) \quad p_{i1} = \frac{e^{X'_{i1}\beta}}{1 + e^{X'_{i1}\beta} + e^{X'_{i2}\beta}}$$

$$(2.40) \quad p_{i0} = \frac{1}{1 + e^{X'_{i1}\beta} + e^{X'_{i2}\beta}}$$

A very important result of McFadden (1974), which shows how the unordered independent logit model is derivable from utility maximization. Suppose that a particular individual i whose utilities associated with three alternatives

$$(2.41) \quad U_{ij} = \mu_{ij} + \epsilon_{ij} \\ j = 0, 1, 2 \quad i = 1, 2, \dots, n$$

Where μ_{ij} is a nonstochastic function of explanatory variables and unknown parameters and ϵ_{ij} is an unobservable random variable. McFadden proved that letting ϵ_j for ϵ_{ij} , that the model of the form (2.38)-(2.40) is derived from utility maximization if and only if $\{\epsilon_j\}$ are independent and the distribution function of ϵ_j is given by $\exp(-e^{-\epsilon_j})$. This distribution is called Type I extreme value distribution, or log Weibull distribution, by Johnson and Kotz (1970a, p.272).

Its density is given by $e^{-\epsilon_j} \exp(-e^{-\epsilon_j})$, which has a unique mode at zero and a mean of approximately 0.577.

Denoting this density by $f(\cdot)$,

$$(2.42) \quad p(y_i=2) = p(U_{i2} > U_{i1}, U_{i2} > U_{i0})$$

$$p(y_i=1) = p(U_{i1} > U_{i2}, U_{i1} > U_{i0})$$

$$p(y_i=0) = p(U_{i0} > U_{i2}, U_{i0} > U_{i1})$$

Considering the case of $p(y_i=2) = p(U_{i2} > U_{i1}, U_{i2} > U_{i0})$

$$= p(\epsilon_2 + \mu_2 - \mu_1 > \epsilon_1, \epsilon_2 + \mu_2 - \mu_0 > \epsilon_0)$$

$$= \int_{-\infty}^{\infty} f(\epsilon_2) \left[\int_{-\infty}^{\epsilon_2 + \mu_2 - \mu_1} f(\epsilon_1) d\epsilon_1 \cdot \int_{-\infty}^{\epsilon_2 + \mu_2 - \mu_0} f(\epsilon_0) d\epsilon_0 \right] d\epsilon_2$$

$$= \int_{-\infty}^{\infty} e^{-\epsilon_2} \exp(-e^{-\epsilon_2}) \cdot \exp(-e^{-\epsilon_2 - \mu_2 + \mu_1}) \cdot \exp(-e^{-\epsilon_2 - \mu_2 + \mu_0}) d\epsilon_2$$

$$= \frac{e^{\mu_{i2}}}{e^{\mu_{i0}} + e^{\mu_{i1}} + e^{\mu_{i2}}} = \frac{e^{\mu_{i2} - \mu_{i0}}}{1 + e^{\mu_{i1} - \mu_{i0}} + e^{\mu_{i2} - \mu_{i0}}}$$

$$= \frac{e^{x'_{i2}\beta}}{1 + e^{x'_{i1}\beta} + e^{x'_{i2}\beta}}$$

$$\text{for } \mu_{i1} - \mu_{i0} = x'_{i1}\beta \text{ and } \mu_{i2} - \mu_{i0} = x'_{i2}\beta$$

III SOME APPLICATIONS OF THE QR4

Qualitative Response Models are applied in various fields, some are considered in this chapter.

1. Lee and DAHM [15].

The purpose here is to estimate the amount of Guthion residue present in the test preparation, for comparison, the standard preparation Guthion is made to contain a given amount of control extract in order that the masking effect due to plant lipids and other inactive substances present in the test preparation be the same for both preparations. All doses of the test preparation contain the same total amount of plant extract, one part due to the test extract itself, and one part due to the control extract added. The flies were observed at 17 hours exposure to count the numbers alive, moribund, and dead respectively.

A test using minimum χ^2 estimate that (the hypothesis) the slope of moribund and dead are equal using trichotomous model or pooling moribund with dead using probit transformation (normit) and logit transformation lead to a similar result.

Conclusion

If the response in a biological assay is polychotomous, it is more efficient to use this information explicitly in

analyzing the data rather than pool certain outcomes in order to make the response dichotomous.

2. Hutchens' Empirical Work [17]

Utilizes data on female heads with children drawn from the "Michigan" data and is restricted to 20 states with large AFDC (Aid to Family with Dependent Children) populations. Hypotheses were tested through maximum likelihood estimation of logistic model of the form

$$P = 1 / (1 + e^{-BX})$$

where P is the probability of remarriage over two years, X is a vector of exogeneous variable and B is a vector of estimated coefficients.

The probability of AFDC receipt in 1970 was estimated with a logistic function, The estimated model is

$$\begin{aligned} \log e \left(\frac{\hat{p}}{1-\hat{p}} \right) &= 0.7736 - 0.35616X_1 - 0.00300X_2 + 0.005143X_3 \\ &\quad (1.2213) \quad (2.5355) \quad (2.0806) \quad (2.8172) \\ &\quad + 0.00008X_4 - 0.20154X_5 - 0.03455X_6 \\ &\quad (0.25281) \quad (0.51559) \quad (2.776) \\ &\quad + 0.03715X_7 + 0.8967X_8 \\ &\quad (1.2679) \quad (3.3910) \end{aligned}$$

where

- P = the probability of AFDC receipt
- X_1 = the wage rate
- X_2 = non-wage income
- X_3 = AFDC guarantee
- X_4 = AFDC breakdown
- X_5 = no earnings (binary variable)
- X_6 = age of female head
- X_7 = presence of children under five (binary variable)
- X_8 = Variable including disability (binary variable)

The empirical results showed that AFDC transfers reduce the probability of remarriage.

Concluding Remarks

Theoretically, AFDC transfers should reduce the probability of participation in marital search and increase the duration of search for female heads with children. Both effects will tend to reduce the probability of remarriage over a short time period. The empirical work presented here indicates that an increase in the level of the AFDC does indeed reduce the probability of remarriage over two years.

3. Li (1977) [21]

Consider the distribution between the housing consumption of homeownership and renters. There are a large number of studies focusing on the explanation of housing tenure

status, which probably is the most important decision about the nature of housing consumption. Income, family size, age of head, and race of head are generally found to be primarily determinants of homeownership.

Previous studies, however, have employed a linearly additive regression model having a dichotomous (0 and 1) dependent variable, which is inconsistent with the expectation of non-linear effects and interaction effects because the probability is bound between 0 and 1. Existing theory of housing consumption also suggests that certain interaction effects should not be ignored. It is, therefore, imperative to test if significant interactions exist.

- i) between age of head and family size as suggested by the life-cycle hypothesis;
- ii) between family size and income resulting from the budget constraint;
- iii) between age of head and income because their joint effect may serve as a proxy for wealth, and
- iv) between race and income because of both the income effect and the substitution on the consumption of black households that face a higher relative price for housing as a result of racial discrimination.

The author adopted the logit model for the analysis of homeownership. The Berkson-Theil method, which employs the cell frequency distribution to derive the logit estimate, is asymptotically equivalent to the maximum likelihood procedure on computational ground. By a monotonic

transformation of a probability having finite range (0,1) to a logit having infinite range $(-\infty, \infty)$, the problem of heteroscedasticity in the error term associated with a regression having a dichotomous dependent variable is avoided. Specifically, the logit is defined as the natural logarithmic value of the odds in favor of a positive response, that is

$$L_i = \log \frac{p_i}{1-p_i} = \beta_0 + \beta X_i$$

where p_i is conditional probability of a positive response with characteristic X_i , and β 's are parameters. It is easily seen that

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta X_i)}}$$

Since the true logit L_i is not observed,

$$\hat{L}_i = \log \frac{f_i}{1-f_i} = \log \frac{p_i}{1-p_i} + U_i$$

where f_i = the observed relative frequency in cell i .

$$U_i = \hat{L}_i - L_i \text{ is the error term}$$

where f_i asymptotically $N(p_i, p_i(1-p_i))$

$$E(U_i) = 0, V(U_i) = \frac{f_i(1-f_i)}{n_i}$$

The generalized least squares (GLS) estimators of β are, in matrix notation, given by

$$b = (X'V^{-1}X)^{-1} X'V^{-1}\hat{L} \quad (*)$$

V = diagonal covariance matrix of the error term.

A standard chi-square test for testing the validity of the logit specification can be obtained from the discrepancies b/n the observed relative frequencies and the estimated probabilities. That is for large sample based on equation (*),

$$\chi^2 = (\hat{L} - Xb)'V^{-1}(\hat{L} - Xb)$$

with degrees of freedom equals to N-k where

N = number of cells

k = number of parameters

An additive logit model of homeownership: The logit model for testing non linear effects and interaction effects specifies that the natural logarithm of the odds in favor of homeownership is a function of income, age of head, family size, and race of head. Specifically, the additive model can be expressed as:

$$L_i = \log \frac{f_i}{1-f_i} = \beta_0 X_0 + \sum_j \beta_{1j} X_{1ji} + \sum_k \beta_{2k} X_{2ki} \\ + \sum_m \beta_{3m} X_{3mi} + \sum_n \beta_{4ni} X_{4ni} + U_i \quad (**)$$

where f_i denotes the cell relative frequency that household type i is a homeowner;

X_0 denotes the constant term, i.e. age of head under 25, income less than \$5,000, two persons, white husband-wife family;

X_{1ji} is a set of four dummy variables denoting five categories for age of head under 25, 25-34, 35-44, 45-64 and over 65;

X_{2ki} is a set of three dummy variables denoting four income class: less than \$5,000, \$5,000-9,999, \$10,000-14,999 and over \$15,000;

X_{3mi} is a set of three dummy variables denoting four family sizes: 2 persons, 3-4 persons, 5 persons, and 6 or more persons;

X_{4ni} denotes race of head: $X_{42}=0$ if nonblack; $X_{42}=1$ if black.

Equation (**) tested empirically with data for husband-wife families in Boston SMSA and in the Baltimore SMSA available from 1970 census metropolitan housing characteristics. Boston and Baltimore are chosen because of the contrast in the relative size of the black population and the difference in the rate of homeownership for the two racial groups. Baltimore SMSA is about four-fifths the size of Boston SMSA, yet, has nearly four times as many blacks. The rates of homeownership are 23.6 and 35.5 percent for Blacks in Boston and Baltimore, respectively;

where as the rates of homeownership for the whites are 60.8 and 69.5 percent respectively.

Estimation of logit coefficients of homeownership for husband-wife families for the basic additive model (***) is given on page 1082 (Table 1) [21] and the detail is given.

Summary

The paper examines two assumptions about the logit specification of homeownership: the assumption of linear effects and additive effects of the three variables (age of head, income and family size), the assumption of a linear effect is statistically most rejectable, where as, the assumption of a linear income effect is less rejectable. The former increases the goodness of fit chi-square value above the basic additive model by about 50%. Furthermore, in two of the three linear-additive models (Age+income+size + race, Age + income + size + race and Age + income + size + race), the assumption of two linear effects in a doubling of the chi-square value for the same additive model.

Of the size two dimensional interactions, the income-size interaction and the age-size interaction are statistically most significant. The former is to be expected from budget constraint; the latter is consistent with the family life-cycle hypothesis. The inclusion of either interaction

into the basic additive model reduces the chi-square value by more than 20 percent. It is not surprising, therefore, to find the allowance for the age-income-size three dimensional interaction reduces the chi-square value from that for the basic model by 96% in the Boston estimate and by 86 percent in the Baltimore estimate. On the other hand, the interactions between the race dummy variables and other explanatory variables, attributable to racial discrimination, are found to be secondary importance relative to the income-size and age-size interactions. The importance of such interactions, however, seems to increase as the number of black families in a SMSA expands as suggested by the difference between the estimates of Boston and the estimates of Baltimore.

II. DATA

The method of data collection is a multistage stratified cluster sampling and the unit of observation is a household [20]. The country is classified into south and North so as to capture distinct ethnic, cultural, religious and economic variables. There are nine administrative regions in the south and five in the North and one region each was randomly selected. Within each region there are subregions and one subregion was again selected; within each sub regions there are villages and within each village one gets peasant associations. Five peasant association from south and four from the North were selected. Between forty and fifty percent of the members of peasant associations were selected for interview. In the end we had 843 observations from the south and 801 from the North.

The final result showed a clear distinction between the two regions. Peasants in the south are cash crop cultivators while those in the North are subsistence crop farmers. The land in the North is highly fragmented while peasants in the south have their land in one plot. Household members in the south are primarily protestants (converted from paganism by Norwegian Missionaries more than 50 years ago). Those in the North belong to orthodox christianity. There are a number of Moslem groups in both regions. Since in the south cash crop producers do not use fertilizers, in this work observations from the North region only are considered.

V. APPLICATION

5.1 Utility Maximization Theory

The study will be based on a standard microeconomic utility maximization theory. It will be assumed that the objective of a peasant household will be to maximize his utility. This has been applied by Quandt (1975), Hansman and Wise (1978), Domencich and Mcfadden (1975) and others [20].

Once we develop a stochastic utility regarding a farmer's behaviour we apply the appropriate qualitative response models to see the determinants of fertilizer use in Ethiopia. We will then test the model using socio-economic and demographic variables.

The utility of a farmer with regards to fertilizer usage will be indirectly measured along the ideas developed by Domencich and Mcfadden (1975). We will let U_{i1} be the indirect utility of the i^{th} farmer from fertilizer use while U_{i2} will be the indirect utility from not using fertilizer. Thus

$$U_{i1} = \alpha_0 + X_{i1}'\beta_1 + y_i'\alpha_2 + \epsilon_{i1} \quad (5.1)$$

$$U_{i2} = \beta_0 + X_{i2}'\beta_1 + y_i'\theta_2 + \epsilon_{i2}$$

Where X_i is a vector of a household demographic characteristics, y_i is a vector of household endowments and ϵ_i 's are additive error terms. The introduction of different

constant terms α_0 and β_0 for each of the utility function is to capture the effect of variables not included in the X_i and Y_i vectors.

The farmer will use fertilizer if $U_{i1} > U_{i2}$ and will not use if $U_{i2} > U_{i1}$. We will assume that the probability of getting the same utility from the two choices as being zero. In other words, $P(U_{i1} = U_{i2}) = 0$. We now define

$$Y_i = \begin{cases} 0 & \text{if fertilizer is not used} \\ 1 & \text{if fertilizer is used} \end{cases} \quad (5.2)$$

$$\begin{aligned} P\{Y_i=1\} &= P\{U_{i1} > U_{i2}\} \\ &= P\{\epsilon_{i2} - \epsilon_{i1} < \alpha_0 - \beta_0 + (X_{i1} - X_{i2})' \beta_1 \\ &\quad + Y_i' (\alpha_2 - \beta_2)\} \\ &= F(\alpha_0 - \beta_0 + (X_{i1} - X_{i2})' \beta_1 + Y_i' (\alpha_2 - \beta_2)) \end{aligned} \quad (5.3)$$

Where F is the distribution function of $\epsilon_{i2} - \epsilon_{i1}$. Equation (5.3) implies that the choice of a particular qualitative response model is similar to the assumed distribution function of the difference between the two error terms.

The above utility models (as stated in Chapter 2) can be indirectly estimated using qualitative response models, such as the linear probability, probit, and logit models; i.e. if the distribution of $\epsilon_{i2} - \epsilon_{i1}$ assumed to be normal, logistic or uniform on $(0,1)$ then the response model will be probit, logit and linear probability models respectively.

Since peasants in the south are cash crop cultivators, no use of fertilizers. So for this research work data from the North region is considered.

A descriptive statistics for selected variables is given in Table 1.

Table 1:

Variable	Definition	Obs	Mean	Std. Dev.
tarea	total area	723	204.57	194.39
cattle	No. of cattle	800	3.26	3.38
heduc	education of head	748	1.45	.66
equip'n	No. of equipments	800	724.00	4.83
tm1014	total No. of males aged 10-14	800	0.34	0.60
tm1554	total No. of females aged 15-54	800	1.35	1.17
plots	No. of plots	712	3.63	2.15
nfincome	annual nonfarm income (birr)	788	23.04	113.23
fertilizer	amount of fertilizer used in kg.	801	35.87	62.54

We see that the number of observations vary for each variables, which is due to non-response and incompleteness of the questionnaire. When the data processed by the computer only 661 observations are considered.

8.2 Variables

The following independent variables are considered in the model.

Demographic Variables:

	<u>Notations</u>
i) Education level of household head	heduc
ii) Total number of females in the household aged 10-14	tf1014
iii) Total number of males in the household aged 10-14	tm1014
iv) Total number of males in the household aged 15-54	tm1554
v) Total number of females in the household aged 15-54	tf1554

Endowments:

i) Total farm area of the household	tarea
ii) Cattle	cattle
iii) Number of farm equipments	equipn
iv) Number of plots of land, and Non-farm income	plots nfincome

where

$$\text{heduc} = \begin{cases} 1 & \text{if the head is illiterate} \\ 2 & \text{if the head reads and writes} \\ 3 & \text{if the head grade 1-4} \\ 4 & \text{if the head is over grade 4} \end{cases}$$

- nfincome = amount of non-farm income in birr (annually)
- plots = number of plots of land which indicates the level of fragmentation of land
- tarea = estimated total size of a household's land in timad*
- equipn. = number of equipments (such as yoke, plough-share, pade, etc.), which is mostly used as an indicator of a household's wealth.

5.3 Estimation Methods and Results

The data from this region is processed using a program (statistical software) known as STATA. We let X_i to be a vector of explanatory variables, the result shows that about 56% of households use fertilizers.

Denoting X_1 , X_2 and X_3 to be vectors of minimum, mean and maximum values of the observations respectively, we obtain the following values

$$X_1' = [1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0]$$

$$X_2' = [1, 204.451, 3.552, 1.416, 8.110, 0.340, 0.355, 1.553, 1.393, 3.641, 13.341]$$

$$X_3' = [1, 1800, 47, 4, 32, 4, 3, 9, 7, 9, 1320]$$

$$\text{for } X_i' = [1, \text{tarea}, \text{cattle}, \text{heduc}, \text{equipn.}, \text{tf1014}, \text{tm1014}, \text{tm1554}, \text{tf1554}, \text{plots}, \text{nfincome}]$$

$$i = 1, 2, \dots, 661$$

*4 timad = 1 hectare.

Applying the three commonly used methods of qualitative response models, namely: Linear probability, probit and logit models, all found to have positive relation with the probability of fertilizer usage, except nfincome. Tables of all the results needed for our work is described below for the three models.

. Linear probability model estimates

Table 2: Determinants of fertilizer use - North-West

Variable	Coefficient	Standard Error	t
tarea	0.00009	0.00009	0.972
Cattle	0.01332	0.00557	2.392
heduc	0.00914	0.02870	0.318
equipment	0.03190	0.00455	7.016
tf1014	0.00083	0.02940	0.028
tm1014	0.04263	0.02916	1.462
tm1554	0.01888	0.01615	1.169
tf1554	0.02402	0.01584	1.517
plots	0.03630	0.00865	4.198
nfincome	-0.00025	0.00019	-1.269
constant	0.01459	0.06210	0.235

n = 661

F(10, 650) = 16.51

R² = 0.2025

Root MSE = 0.44704

Prob > F = 0.000

adj. R² = 0.1903

Table 3: Logit Model Estimates

Variable	Coefficient	Standard Error	t
tarea	0.00048	0.00050	0.961
cattle	0.10295	0.03755	2.742
heduc	0.05634	0.15053	0.374
equipment	0.17197	0.0264	6.520
tf1014	0.02465	0.15191	0.162
tml014	0.20093	0.15135	1.328
tml554	0.11061	0.08630	1.282
tf1554	0.11266	0.08149	1.383
plots	0.17228	0.0442	3.898
nfincome	-0.00155	0.00131	-1.189
constant	-2.6220	0.35445	-7.397

n = 661

chi 2(10) = 155.75

prob > chi 2 = 0.000

log likelihood = 375.56136

We can observe that for the probit and logit models the Chi-squares and log likelihood values are almost equal.

Let us denote now the vector of coefficients of the linear probability, probit and logit model by β_{lp} , β_{pr} and β_{log} respectively. Where

β_{lp} = contains column 2 of table 1 taking the constant at the top.

Table 4: Probit Model Estimates

Variable	Coefficient	Standard Error	t
tarea	0.00030	0.00030	1.009
cattle	0.05603	0.02158	2.596
heduc	0.04059	0.09040	0.449
equipment	0.10281	0.01537	6.687
tf1014	0.02026	0.08955	0.227
tm1014	0.10355	0.08961	1.156
tm1554	0.06642	0.05148	1.290
tf1554	0.07137	0.04879	1.463
plots	0.10439	0.02625	3.976
equipment	-0.00097	0.00079	-1.216
constant	-1.576028	0.20733	-7.602

n = 661
Chi 2(10) = 154.94
prob Chi 2 = 0.000
log likelihood = -375.96905

We can observe that for the probit and logit models the Chi-square and log likelihood values are almost equal.

Let us denote now the vector of coefficients of the linear probability, probit and logit model by β_{Lp} , β_p and β_L respectively. Where

β_{Lp} = contains column 2 of table 2 taking the constant at the top.

β_L = contains column 2 of Table 3 taking
the constant at the top.

β_{LP} = contains column 2 of Table 4 taking
the constant at the top.

5.4 Discussion

Demographic Variables

Households with more dependents, more adults and in general large household have a high probability of using fertilizers than household with less dependents and less number of adults. The reason may be that large size households have been established for larger periods and may accuire more land and which in turn may require more fertilizers. Besides large size households have more chance of getting fertilizers on credit than small size households. The latter are usually headed by younger individuals. That is, those who have not^{vet} established their credit worthiness.

Endowments

Households with high amount of capital inputs as measured by the number of farm implements, cattle ownership, total farm area, and number of plots of land tend to have positive effect on fertilizer usage. The above variables are indicators of how much an Ethiopian farmer

is devoted in the process of agricultural productivity. In other words, the higher the endowments, the more inputs in the process of agricultural production and may in turn implies higher probability of fertilizer use. Therefore, endowments are indicators of a household's potential to use fertilizers.

On the other hand individuals with more non-farm income (such as traders, black smiths, daily labourers) have other sources of earnings eventhough they are basically farmers. Since they have other sources of income, chances are that they are less likely to engage and invest on agriculture. The results are therefore consistent with the hypothesis. In otherwords people who have non-farm income are less likely to use fertilizers.

5.5 Probability Generation

Considering these models we will try to generate the probability of fertilizer usage, using minimum, mean and maximum values of the explanatory variables. For

$$y_i = \begin{cases} 0 & \text{if a household doesn't use fertilizers} \\ 1 & \text{if a household used fertilizers} \end{cases}$$

$$i = 1, 2, \dots, 661$$

i) Linear probability estimates

Defining

$$P(v_i=1) = \begin{cases} 0 & \text{if } X_i' \hat{\beta}_{LP} \leq 0 \\ X_i' \hat{\beta}_{LP} & \text{if } 0 < X_i' \hat{\beta}_{LP} < 1 \\ 1 & \text{if } X_i' \hat{\beta}_{LP} \geq 1 \end{cases}$$

we obtain,

$$X_1' \hat{\beta}_{LP} = 0.023727$$

$$X_2' \hat{\beta}_{LP} = 0.997382$$

$$X_3' \hat{\beta}_{LP} = 4.9976239,$$

based on our definition the probabilities using minimum mean and maximum values to the nearest 3 decimal digits are 0.024, 0.997 and 1.000 respectively.

The vector of coefficients β_D and β_L are obtained using the Newton-Raphson's method of iteration to find the maximum-likelihood estimates. In both cases the 4th iterations deliver the desired results.

ii) Probit model estimates

$$P(v_i=1) = \Phi(X_i' \hat{\beta}_P)$$

Where Φ stands for the cumulative standard normal distribution, for the values

$$X_1' \hat{\beta}_P = -1.47322$$

$$X_2' \hat{\beta}_P = -0.15561$$

$$X_3' \hat{\beta}_P = 4.23357$$

Then the probabilities that a farmer will use fertilizers are 0.0718 using minimum values, 0.4336 using mean values and approximately 1 when the maximum values are considered.

iii) Logit model estimates

$$P(y_i = 1) = \frac{e^{X_i' \hat{\beta}_L}}{1 + e^{X_i' \hat{\beta}_L}}$$

Where for the minimum, mean and maximum values

$$X_1' \hat{\beta}_L = -2.5657$$

$$X_2' \hat{\beta}_L = 0.35123$$

$$X_3' \hat{\beta}_L = 10.79864$$

$$\text{and } \log(p_i / (1 - p_i)) = X_i' \hat{\beta}_L$$

In this case the probabilities for minimum, mean and maximum values that a farmer will use fertilizers are 0.0713, 0.5869 and 0.9999 respectively.

We see that probit and logit models give similar results whereas the linear probability differ much from both.

Remarks: Minimum, mean and maximum probability values

Once the model is estimated we tried conditions under which the probability of using fertilizer is minimum and maximum. The minimum probability for probit and logit is

about the same; that is a value of 0.072 and 0.071. This is generated when identifying farmers whose land size, number of cattle, farm equipments, number of dependents, plots, non-farm income as well as the level of education of the household head is at its minimum.

The maximum probability for probit and logit models is again almost the same and is about unity.

Simulation methods

The results obtained above may help policy makers to identify farmers that are likely to be users of fertilizers. They may also try to give incentives to non-users so as to use fertilizers and thereby increase agricultural yield. The policy options are nothing but the explanatory variables already mentioned. It may not be feasible to change all the explanatory variables simultaneously.

Below we vary some of the feasible explanatory variables and try to generate the probability of using fertilizers. The variables that we will vary are

1. Education of household head (heduc)
2. Total area (tarea)
3. Non-farm income (nfincome)
4. Cattle
5. Farm equipment (equirn)

Note that demographic variables have been omitted from simulation because the results would require higher fertility so as to use fertilizers.

The simulation is done for four different cases, where the mean values of heduc, tarea, nfincome, cattle and equipm increased by 0.5 sd, 1.0 sd, 1.5 sd and 2.0 sd (sd stands for standard deviation of the variable), and when the demographic variables and plots remain fixed at their mean values.

Table 5

Rate of Increase (R)	heduc (mean+R)	tarea (mean+R)	nfincome (mean+R)	cattle (mean+R)	equipm (mean+R)	Over all prob. Probit	logit
0.5 sd	1.781	301.765	79.66	4.952	9.655	0.6554	0.667
1.0 sd	2.111	398.94	136.275	6.644	12.070	0.7688	0.781
1.5 sd	2.411	496.155	192.89	8.336	14.485	0.8577	0.871
2.0 sd	2.771	5.9335	249.505	10.08	16.908	0.9207	0.918

The probit and logit models generate almost the same probabilities in all cases. On the other hand we observed that if farmers have more land, cattle, equipment, non-farm income and higher level of education the probability (the potential) that a farmer will use fertilizer will increase as well. That is the higher the rate of increase of this explanatory variables the higher the probability of a household's fertilizer usage.

5.6 Comparison of Models

To compare the probit and logit models, we consider here only the log-likelihood function

$$\ell = \sum_{i=1}^n [y_i \log \hat{F}_i + (1-y_i) \log(1-\hat{F}_i)]$$

where $\hat{F}_i = \pi(X_i' \hat{\beta}_{ML})$, $i=1,2,\dots,661$ and $\hat{\beta}_{ML}$ denotes the maximum likelihood estimates. This measure has an intuitive appeal [2].

The log likelihood values for the probit and logit models are -375.96905 and -375.56136 respectively. This strongly confirms that the probit and logit models give similar results. On the other hand, as discussed in Chapter 2, the linear probability model estimates are different from the two models and the coefficient of determination (R^2) and adjusted R^2 have values 0.2025 and 0.1903 respectively, which indicates that the explanatory power of the model for the given data is not strong. In empirical results it is not uncommon to observe low explanatory power for linear regression but high value of χ^2 for probit and logit. Thus results are not internally inconsistent.

5.7 Tests of Hypotheses

Two types of hypothesis are tried for the estimated logit and probit models. First each of the coefficients was tested using standard normal distribution. Second

the overall estimated equation was tested using standard χ^2 test.

To test the hypothesis that each b_i (coefficient of the explanatory variable) equals to zero,

the necessary $H_0: b_i = 0$ and $H_A: b_i \neq 0$ are stated

avoid the term $H_A: b_i \neq 0, i=1,2,\dots,11$ for both

where b_1 stands for the constant term and 11 is the number of parameters (b_i), using Z test

$$Z = \hat{b}_i / Sb_i, \alpha = 0.05$$

Where Sb_i = standard error of b_i , and for each model calculated values are listed under the column of t in Tables 3 and 4.

To test the hypothesis

$$H_0: Q'\beta = 0 \quad (\text{the explanatory variables have no effect on the probability of fertilizer usage})$$

$$H_A: Q'\beta \neq 0$$

where

$$Q' = \begin{bmatrix} 0' \\ I \end{bmatrix}$$

$$Q = 0_{11 \times 1} \quad \text{a vector of zeros}$$

$$I = \text{an identity matrix of rank 10.}$$

The null hypothesis states that all the b_i 's are zero except the constant term. Using Wald's statistic:

$$\begin{aligned} \text{Wald} &= (O'\hat{\beta} - C)'(O'V\hat{\beta}O)^{-1}(O\hat{\beta} - C) \\ &= (O'\hat{\beta})'(O'V\hat{\beta}O)^{-1}(O\hat{\beta}) \end{aligned}$$

Where $V\hat{\beta}$ is the asymptotic variance - covariance matrix of the vector of estimated coefficients (see Appendix A). All the necessary adjustments and rearrangements are done to avoid the terms containing the constant term from both the vector of coefficients $\hat{\beta}$ and variance - covariance matrix $V\hat{\beta}$.

i) Probit model

$$H_0: Q'\beta_n = 0$$

$$H_A: Q'\beta_n \neq 0$$

$$(O'\hat{\beta}_n)'(O'V\hat{\beta}_nO)^{-1}(O\hat{\beta}_n) = 154.94$$

where $\hat{\beta}_n$ = estimated vector of coefficients of the probit model.

$V\hat{\beta}_n$ = asymptotic variance covariance matrix of $\hat{\beta}_n$.

ii) Logit model

$$H_0: Q'\beta_L = 0$$

$$H_A: Q'\beta_L \neq 0$$

$$(O'\hat{\beta}_L)'(O'V\hat{\beta}_LO)^{-1}(O\hat{\beta}_L) = 155.75$$

where $\hat{\beta}_L$ = estimated vector of coefficients for the logit model.

$V\hat{\beta}_L$ = asymptotic variance - covariance matrix of $\hat{\beta}_L$.

From (i) and (ii) we see that at any level of significance (α) the calculated value found to be greater than χ^2_{10} . Therefore we reject the null hypothesis H_0 in both cases, i.e. the probability that a farmer (a household) will use fertilizers depends on the total size of farm land, number of cattle, level of education of the household head, number of farm implements, and family size (children, adults) positively and with a household's non-farm income negatively.

The socio-demographic variables are level of education of household head, number of farm implements, number of cattle, number of plots, farm land size and non-farm income. The following paragraphs summarized the findings.

The explanatory variables satisfactorily explain a farmer's fertilizer usage and this is supported by the test, although there may be some variables which are not included here. First non-farm income all variables seem to affect the probability positively, whereas the farmer had a negative effect on fertilizer usage. This indicates that if a household has earnings from non-farm employment, most of a working householdly farmer, he has no work devoted to agriculture in several seasons. Most of the farmers are to increase their agriculture use and to modernize. The farmer may have more implements including tools to use as he has increased the level of productivity. The agricultural use may not be profitable.

CONCLUSIONS

By hypothesizing that a farmer's demand to use fertilizers is dependent on a household's demographic and non-demographic (endowments) variables, we applied the concept of Qualitative Response Models to study the hypothesis. The demographic variables are number of females aged 10-14, number of males aged 10-14, number of males aged 15-54 and number of females aged 15-54 in a household.

The non-demographic variables are level of education of the household's head, number of farm implements, number of cattle, number of plots, farm land size and non-farm income. The following paragraphs summarize the findings.

The explanatory variables satisfactorily explain a farmer's fertilizer usage and this is supported by the test, even though there may be some variables which are not included here. Except non-farm income all variables seem to affect the probability positively, whereas the former has a negative effect on fertilizer usage. This indicates that if a household has earnings from non-farm employment, even if a person is basically a farmer, he has not much devotion to agriculture due to several reasons. Some of the reasons may be

- i) income from agriculture may not be satisfactory;
- ii) the farmer may lack farm implements including oxen so that he has to rent his land on contract basis;
- iii) the household's farm land may not be productive;

iv) agricultural production in the region may be altered due to man made and natural problems.

Using the minimum, mean and maximum values of the explanatory variables we have tried to generate probabilities. If a farmer has the mean values of the mentioned explanatory variables, the probabilities found to be 0.44 in the case of probit and 0.58 in the case of logit model which is around 0.5. This would mean that a farmer is less likely to use fertilizers. On the other hand using the maximum value of the explanatory variables, the probabilities are around unity. This is true for the three models. However, taking into account the case of population pressure in Ethiopia, it may not be worthwhile to suggest an increase in a household's family size.

We also tried to generate probabilities for various values of the explanatory variables by keeping demographic variables and plots at mean values. This leads to an increase in probabilities as shown in Table 5. The case of non-farm income seems to give a different result; however, we can conclude that if a farmer has these incentives the last variable has little to do with the probability of fertilizer usage.

Finally, since increase in the use of fertilizers is directly related to increase in agricultural production and which in turn implies a solution to be self sufficient

APPENDIX A: Table of logit model estimates

. logit y tarea cattle heduc equipn tf1014 tm1014 tm1554 tf1554 plots nfin

Iteration 0: Log Likelihood = -453.43811
Iteration 1: Log Likelihood = -379.97588
Iteration 2: Log Likelihood = -375.65993
Iteration 3: Log Likelihood = -375.56144
Iteration 4: Log Likelihood = -375.56136

Logit Estimates

Number of obs = 661
chi2(10) = 155.75
Prob > chi2 = 0.0000

Log Likelihood = -375.56136

Table with 6 columns: Variable, Coefficient, Std. Error, t, Prob > |t|, Mean. Rows include y, tarea, cattle, heduc, equipn, tf1014, tm1014, tm1554, tf1554, plots, nfincome, and _cons.

. correlate, _coef covariance

Table showing covariance matrix for variables: tarea, cattle, heduc, equipn, tf1014, tm1014, tm1554, tf1554, plots, nfincome, and _cons.

APPENDIX B: Table of probit model estimates

. probit y tarea cattle heduc equipn tf1014 tm1014 tm1554 tf1554 plots nfin

Iteration 0: Log Likelihood = -453.43811
 Iteration 1: Log Likelihood = -379.31799
 Iteration 2: Log Likelihood = -375.99682
 Iteration 3: Log Likelihood = -375.96906
 Iteration 4: Log Likelihood = -375.96905

Probit Estimates

Number of obs = 661
 chi2(10) = 154.94
 Prob > chi2 = 0.0000

Log Likelihood = -375.96905

Variable	Coefficient	Std. Error	t	Prob > t	Mean
y					.5597579
tarea	.0003064	.0003038	1.009	0.314	204.4508
cattle	.056031	.0215829	2.596	0.010	3.552194
heduc	.0405927	.0904004	0.449	0.654	1.416036
equipn	.1028105	.0153747	6.687	0.000	8.110439
tf1014	.0202649	.0893538	0.227	0.821	.3403933
tm1014	.1035584	.0896121	1.156	0.248	.3555219
tm1554	.0664252	.0514888	1.290	0.197	1.553707
tf1554	.0713783	.0487963	1.463	0.144	1.393343
plots	.1043951	.0262593	3.976	0.000	3.641452
nfincome	-.0009718	.0007989	-1.216	0.224	13.34144
_cons	-1.576028	.2073301	-7.602	0.000	1

. corr, _coef covariance

	tarea	cattle	heduc	equipn	tf1014	tm1014	tm1554
tarea	9.2e-08						
cattle	-6.6e-07	.000466					
heduc	-1.5e-06	-.000146	.008172				
equipn	-4.4e-07	-.000089	-.000049	.000236			
tf1014	-7.1e-07	.000045	.000467	-.000011	.007984		
tm1014	5.4e-07	-.000079	.000054	-.000117	.000298	.00803	
tm1554	-3.2e-07	-.000062	.000124	-.000036	-.000436	-.000083	.002651
tf1554	-1.3e-06	.000028	.000348	.000019	-.000431	-.000616	-.000713
plots	-1.3e-06	.000017	.000021	-.000059	.0001	.000015	-.000132
nfincome	-3.3e-08	-8.2e-07	-1.1e-06	1.8e-06	-2.6e-06	-1.3e-06	-2.1e-06
_cons	-3.0e-06	-.000498	-.011143	-.001141	-.002442	-.00097	-.001936

	tf1554	plots	nfincome	_cons
tf1554	.002381			
plots	.000035	.00069		
nfincome	1.0e-06	3.4e-07	6.4e-07	
_cons	-.002433	-.001702	-8.1e-06	.042986

APPENDIX C: Table of linear probability

56

	Model estimates				
tm1014	.2009337	.1513563	1.328	0.185	.3555219
tm1554	.110612	.0863016	1.282	0.200	1.553707
tf1554	.1126659	.0814943	1.383	0.167	1.393343
plots	.172284	.0442	3.898	0.000	3.641452
nfincome	-.0015563	.001309	-1.189	0.235	13.34144
_cons	-2.622016	.3544548	-7.397	0.000	1

. reg y tarea cattle heduc equip tf1014 tm1014 tm1554 tf1554 plots nfinco
(obs=661)

Source	SS	df	MS	Number of obs = 661	
Model	32.990295	10	3.2990295	F(10, 650) =	16.51
Residual	129.899266	650	.199845025	Prob > F =	0.0000
Total	162.889561	660	.246802366	R-square =	0.2025
				Adj R-square =	0.1903
				Root MSE =	.44704

Variable	Coefficient	Std. Error	t	Prob > t	Mean
y					.5597579
tarea	.0000932	.000096	0.972	0.332	204.4508
cattle	.0133185	.0055687	2.392	0.017	3.552194
heduc	.0091403	.0287009	0.318	0.750	1.416036
equipn	.0319031	.0045471	7.016	0.000	8.110439
tf1014	.0008308	.0294025	0.028	0.977	.3403933
tm1014	.0426258	.0291611	1.462	0.144	.3555219
tm1554	.0188775	.0161468	1.169	0.243	1.553707
tf1554	.0240222	.0158395	1.517	0.130	1.393343
plots	.0363017	.0086484	4.198	0.000	3.641452
nfincome	-.0002488	.000196	-1.269	0.205	13.34144
_cons	.0145863	.062103	0.235	0.814	1

REFERENCES

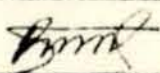
1. Akin, J.S.; Guilkey, B.K. and Sickles, P. "A Random Coefficient Probit Model with an Application to a study of Migration." *J. Econometrics*, Oct.-Dec. 1979, pp. 233-45.
2. Amemiya, T. "Qualitative Response Models." *J. Economic Literatures*, Dec. 1981, pp. 1483-1536.
3. Bartel, A. "The Migration Decision: What role does job mobility play?" *Amer. Econ. Rev.*, Dec. 1979, pp. 775-86.
4. Berkson, J. "Why I prefer Logits to Probits." *Biometrics*, Dec. 1951, pp. 327-59.
5. Boskir, M.J. "A Conditional Logit Model of Occupational Choice." *J. Polit. Econ.*, March-April 1974, pp. 389-98.
6. Cox, B.R. "Some Procedures Connected with the Logistic Qualitative Response Curve." In research papers in statistics. Edited by F.N. David. New York: Wiley, 1966, pp. 55-71.
7. Crag, J.G. and Uhler, R.S. "The Demand for Automobiles." *Can. J. Econ.*, Aug. 1970, pp. 386-406.
8. Daganzo, C. *Multinomial Probit*. New York: Academic Press, 1979.
9. Deacon, R. and Shviro, P. "Private Preference for Collective Goods Revealed through Voting and Referenda." *Amer. Econ. Rev.*, Dec. 1975, pp. 943-55.
10. Dobson, A.G. *Introduction to Statistical Modelling*. New York: Chapman and Hall, 1983.
11. Debreu, G. "Review of 'Individual Choice Behavior' by R. Luce." *Amer. Econ. Rev.*, March 1960, pp. 186-88.
12. Domencich, T.A. and McFadden, D. *Urban Travel Demand*. Amsterdam: North Holland, 1975.
13. Finney, D.J. *Probit Analysis*. Third edition, Cambridge: University Press, 1971.

14. Goldberger, A.S. "Correlations Between Binary Outcomes and Probabilistic Predictions." J. Amer. Statist. Assoc., March 1973, pp. 84.
15. Gurland, J.; Lee, I. and Bahr, P.A. "Polychotomous Quantal Response in Biological Assay." Biometrics, Sept. 1960, pp. 382-98.
16. Haberman, S.J. Analysis of Qualitative Data, Vol. 11, new developments. New York: Academic Press, 1979.
17. Hutehens, R.M. "Welfare, Remarriage and Marital Search." Amer. Econ. Rev., June 1979, pp. 369-79.
18. Johnston, J. Econometric Methods. Third edition. Singapore: McGraw-Hill.
19. Kidane, A. (1989). "Differences in Educational Opportunities between Males and Females in Rural Ethiopia." A paper presented at the Conference on "The Family Gender Differences and Development" at Yale University, New Haven, Conn. Sept. 4-6, 1989.
20. _____ . "Determinants of Multiple Cropping in Ethiopia, an Application of Qualitative Response Models." Manuscript, Department of Statistics, Addis Ababa University.
21. Li, M.M. "Alogit Model of Homeownership." Economica, July 1977, pp. 1081-89.
22. Morimune, K. "Comparison of Normal and Logistic Models in the Bivariate Dichotomous Analysis." Econometrica, July 1979, pp. 1957-76.

DECLARATION

I, the undersigned, declare that this thesis is my original work and has not been presented for a degree in any other University.

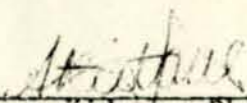
Name Desta W. Mariam

Signature 

Place: Statistics Department

Date: June, 1990

This thesis has been submitted for examination with my approval as University advisor.


Asmerom Kidane, Ph.D.