



Nonparametric identification of regression models containing a misclassified dichotomous regressor without instruments

Xiaohong Chen ^{a,1}, Yingyao Hu ^{b,2}, Arthur Lewbel ^{c,*}

^a Department of Economics, Yale University, Box 208281, New Haven, CT 06520-8281, USA

^b Department of Economics, Johns Hopkins University, 440 Mergenthaler Hall, 3400 N. Charles Street, Baltimore, MD 21218, USA

^c Department of Economics, Boston College, 140 Commonwealth Avenue, Chestnut Hill, MA 02467, USA

ARTICLE INFO

Article history:

Received 6 December 2007

Received in revised form 28 February 2008

Accepted 6 March 2008

Available online 20 March 2008

Keywords:

Misclassification error

Identification

Nonparametric regression

JEL classification:

C14

C20

ABSTRACT

We observe a dependent variable and some regressors, including a mismeasured binary regressor. We provide identification of the nonparametric regression model containing this misclassified dichotomous regressor. We obtain identification without parameterizations or instruments, by assuming the model error isn't skewed.

© 2008 Elsevier B.V. All rights reserved.

1. Motivation

We provide identification of a nonparametric regression model with a dichotomous regressor subject to misclassification error. The available sample information consists of a dependent variable and a set of regressors, one of which is binary and error-ridden with misclassification error that has unknown distribution. Our identification strategy does not parameterize any regression or distribution functions, and does not require additional sample information such as instrumental variables, repeated measurements, or an auxiliary sample. Our main identifying assumption is that the regression model error has zero conditional third moment. The results include a closed-form solution for the unknown distributions and the regression function.

Dichotomous (binary) variables, such as union status, smoking behavior, and having a college degree or not, are involved in many economic models. Measurement errors in dichotomous variables take the form of misclassification errors, i.e., some observations where the variable is actually a one may be misclassified as a zero, and vice versa. A common source of misclassification errors is self-reporting, where people may have

psychological or economic incentives to misreport dichotomous variables (see Bound et al. (2001) for a survey). Misclassification may also arise from ordinary coding or reporting errors, e.g., Kane et al. (1999) report substantial classification errors in both self-reports and transcript reports of educational attainment. Unlike ordinary mismeasured regressors, misclassified regressors cannot possess the properties of classically mismeasured variables, in particular, classification errors are not independent of the underlying true regressor, and are in general not mean zero.

As with ordinary mismeasured regressors, estimated regressions with a misclassified regressor are inconsistent, and the latent true regression model based just on conditionally mean zero model errors is generally not identified in the presence of a misclassified regressor. To identify the latent model, we must either impose additional assumptions or possess additional sample information. One popular additional assumption is to assume the measurement error distribution belong to some parametric family. Additional sample information often used to obtain identification includes an instrumental variable or a repeated measurement in the same sample, or a secondary sample. See, e.g., Carroll et al. (2006), and Chen et al. (2007) for detailed recent reviews on existing approaches to measurement error problems.

In this note we obtain identification without parameterizing errors and without auxiliary information like instrumental variables, repeated measurements, or a secondary sample. A related result is Chen et al. (2008). We show here that, given some mild regularity conditions, a nonparametric mean regression with a misclassified

* Corresponding author. Tel.: +1 617 522 3678.

E-mail addresses: xiaohong.chen@yale.edu (X. Chen), yhu@jhu.edu (Y. Hu), lewbel@bc.edu (A. Lewbel).

¹ Tel.: +1 203 432 5852.

² Tel.: +1 410 516 7610.

binary regressor is identified (and can be solved in closed form) if the latent regression error has zero conditional third moment, as would be the case if the regression error were symmetric. We also briefly discuss how simple estimators might be constructed based on our identification method.

2. Identification

We are interested in a regression model as follows:

$$Y = m(X^*, W) + \eta, \quad E(\eta|X^*, W) = 0 \tag{2.1}$$

where Y is the dependent variable, $X^* \in X = \{0, 1\}$ is the dichotomous regressor subject to misclassification error, and W is an error-free covariate vector. We are interested in the nonparametric identification of the regression function $m(\cdot)$. The regression error η need not be independent of the regressors X^* and W , so we have conditional density functions

$$f_{Y|X^*, W}(y|X^*, w) = f_{\eta|X^*, W}(y - m(X^*, w)|X^*, w). \tag{2.2}$$

In a random sample, we observe $(X, Y, W) \in X \times Y \times W$, where X is a proxy or a mismeasured version of X^* . We assume

Assumption 2.1. $f_{Y|X^*, W, X}(y|X^*, w, x) = f_{Y|X^*, W}(y|X^*, w)$ for all $(x, X^*, y, w) \in X \times X^* \times Y \times W$.

This assumption implies that the measurement error in X is independent of the dependent variable Y conditional on the true value X^* and the covariate W , and so X is independent of the regression error η conditional on X^* and W . This is analogous to the classical measurement error assumption of having the measurement error independent of the regression model error. This assumption may be problematic in applications where the same individual who provides the source of misclassification by supplying X also helps determine the outcome Y , however, this is a standard assumption in the literature of mismeasured and misclassified regressors. See, e.g., Li (2002), Schennach (2004), Mahajan (2006), Lewbel (2007a) and Hu (2006).

By construction, the relationship between the observed density and the latent ones are as follows:

$$f_{Y|X, W}(y|x, w) = \sum_{X^*} f_{Y|X^*, W, X}(y|X^*, w, x) f_{X^*|X, W}(X^*|x, w) = \sum_{X^*} f_{\eta|X^*, W}(y - m(X^*, w)|X^*, w) f_{X^*|X, W}(X^*|x, w). \tag{2.3}$$

Using the fact that X and X^* are 0–1 dichotomous, define the following simplifying notation: $m_0(w) = m(0, w)$, $m_1(w) = m(1, w)$, $\mu_0(w) = E(Y|X=0, w)$, $\mu_1(w) = E(Y|X=1, w)$, $p(w) = f_{X^*=1|X, W}(1|0, w)$, and $q(w) = f_{X^*=1|X, W}(0|1, w)$. Eq. (2.3) is then equivalent to

$$\left(\frac{f_{Y|X, W}(y|0, w)}{f_{Y|X, W}(y|1, w)} \right) = \left(\frac{1 - p(w)}{q(w)} \quad \frac{p(w)}{1 - q(w)} \right) \left(\frac{f_{\eta|X^*, W}(y - m_0(w)|0, w)}{f_{\eta|X^*, W}(y - m_1(w)|1, w)} \right). \tag{2.4}$$

Since $f_{\eta|X^*, W}$ has zero mean, we obtain

$$\mu_0(w) = (1 - p(w))m_0(w) + p(w)m_1(w) \text{ and } \mu_1(w) = q(w)m_0(w) + (1 - q(w))m_1(w). \tag{2.5}$$

Assume

Assumption 2.2. $m_1(w) \neq m_0(w)$ for all $w \in W$.

This assumption means that X^* has a nonzero effect on the conditional mean of Y , and so is a relevant explanatory variable, given W . We may now solve Eq. (2.5) for $p(w)$ and $q(w)$, yielding

$$p(w) = \frac{\mu_0(w) - m_0(w)}{m_1(w) - m_0(w)} \text{ and } q(w) = \frac{m_1(w) - \mu_1(w)}{m_1(w) - m_0(w)} \tag{2.6}$$

Without loss of generality, we assume,

Assumption 2.3. for all $w \in W$, (i) $\mu_1(w) > \mu_0(w)$; (ii) $p(w) + q(w) < 1$.

Assumption 2.3(i) is not restrictive because one can always redefine X as $1 - X$ if needed. Assumption 2.3(ii) implies that the ordering of $m_1(w)$ and $m_0(w)$ is the same as that of $\mu_1(w)$ and $\mu_0(w)$ because $1 - p(w) - q(w) = \frac{\mu_1(w) - \mu_0(w)}{m_1(w) - m_0(w)}$. The intuition of Assumption 2.3 (ii) is that the total misclassification probability is not too large so that $\mu_1(w) > \mu_0(w)$ implies $m_1(w) > m_0(w)$ (see, e.g., Lewbel, 2007a) for a further discussion of this assumption). In summary, we have

$$m_1(w) \geq \mu_1(w) > \mu_0(w) \geq m_0(w).$$

The condition $p(w) + q(w) \neq 1$ also guarantees that the matrix $\begin{pmatrix} 1 - p(w) & p(w) \\ q(w) & 1 - q(w) \end{pmatrix}$ in Eq. (2.4) is invertible. If we then plug into Eq. (2.4) the expressions for $p(w)$ and $q(w)$ in Eq. (2.6), we obtain for $j=0,1$

$$f_{\eta|X^*, W}(y - m_j(w)|j, w) = \frac{\mu_1(w) - m_j(w)}{\mu_1(w) - \mu_0(w)} f_{Y|X, W}(y|0, w) + \frac{m_j(w) - \mu_0(w)}{\mu_1(w) - \mu_0(w)} f_{Y|X, W}(y|1, w). \tag{2.7}$$

Eq. (2.7) is our vehicle for identification. Given any information about the distribution of the regression error η , Eq. (2.7) provides the link between that information and the unknowns $m_0(w)$ and $m_1(w)$, along with the observable density $f_{Y|X, W}$ and observable conditional means $\mu_0(w)$ and $\mu_1(w)$. The specific assumption about η that we use to obtain identification is this:

Assumption 2.4. $E(\eta^3|X^*, W) = 0$.

A sufficient though much stronger than necessary condition for this assumption to hold is that $f_{\eta|X^*, W}$ be symmetric for each $x^* \in X$ and $w \in W$. Notice that the regression model error η need not be independent of the regressors X^* , W , and in particular our assumptions permit η to have heteroskedasticity of completely unknown form.

Let ϕ denote the characteristic function and

$$\phi_{\eta|X^*=j, W}(t) = \int e^{it\eta} f_{\eta|X^*, W}(\eta|j, w) d\eta$$

$$\phi_{Y|X=j, W}(t) = \int e^{ity} f_{Y|X, W}(y|j, w) dy.$$

Then Eq. (2.7) implies that for any real t

$$\ln \left(e^{itm_j(w)} \phi_{\eta|X^*=j, W}(t) \right) = \ln \left(\frac{\mu_1(w) - m_j(w)}{\mu_1(w) - \mu_0(w)} \phi_{Y|X=0, W}(t) + \frac{m_j(w) - \mu_0(w)}{\mu_1(w) - \mu_0(w)} \phi_{Y|X=1, W}(t) \right). \tag{2.8}$$

Notice that

$$\frac{\partial^3}{\partial t^3} \ln \left(e^{itm_j(w)} \phi_{\eta|X^*=j, W}(t) \right) \Big|_{t=0} = \frac{\partial^3}{\partial t^3} \ln \phi_{\eta|X^*=j, W}(t) \Big|_{t=0} = -iE(\eta^3|X^* = j, W = w).$$

Assumption 2.4 therefore implies that for $j=0,1$

$$G(m_j(w)) = 0, \tag{2.9}$$

where

$$G(z) = i \frac{\partial^3}{\partial t^3} \ln \left(\frac{\mu_1(w) - z}{\mu_1(w) - \mu_0(w)} \phi_{Y|X=0, W}(t) + \frac{z - \mu_0(w)}{\mu_1(w) - \mu_0(w)} \phi_{Y|X=1, W}(t) \right) \Big|_{t=0}.$$

This equation shows that the unknowns $m_0(w)$ and $m_1(w)$ are two roots of the cubic function $G(\cdot)$ in Eq. (2.9). Suppose the three roots of this equation are $r_a(w) \leq r_b(w) \leq r_c(w)$. In fact, we have

$$r_a(w) \leq m_0(w) \leq \mu_0(w) < \mu_1(w) \leq m_1(w) \leq r_c(w),$$

which implies bounds on $m_0(w)$ and $m_1(w)$. To obtain point identification of $m_j(w)$, we need to be able to uniquely define which roots of the cubic function $G(\cdot)$ correspond to $m_0(w)$ and $m_1(w)$. This is provided by the following assumption.

Download English Version:

<https://daneshyari.com/en/article/5061764>

Download Persian Version:

<https://daneshyari.com/article/5061764>

[Daneshyari.com](https://daneshyari.com)