**Endogeneity in Stochastic Frontier Models: Copula Approach without Ooutside ~~instruments~~Instruments**

Kien C. Tran[*]
Department of Economics,
University of Lethbridge,
4401 University Drive W
Lethbridge, Alberta,
T1K 3M4 CANADA, email: kien.tran@uleth.ca

~~and~~

Efthymios G. Tsionas,
Department of Economics,
Athens University of Economics and Business
76 Pattison Street,
10434 Athens, GREECE

### Abstract

This papers considers an alternative estimation procedures for estimating stochastic frontier models with endogenous regressors when no external instruments are available. The approach we propose is based on copula function to directly model the correlation between the endogenous regressors and the composed errors. Estimation of model parameters is done using maximum likelihood. Monte Carlo simulations are used to assess and compare the finite sample performances of the proposed estimation procedures.

*JEL Classification:*  C13, C14.

*Keywords*: Stochastic frontier model, Endogenous regressors, Copula function, Maximum likelihood.

## 1. Introduction

A standard approach to handle endogeneity problem in the stochastic frontier models is to use likelihood based instrumental variable estimation methods, see for example, Kutlu (2010), Tran and Tsionas (2013) and Amsler, Prokhorov and Schmidt (2014). This type of approach relies upon the availability of a set of outside information that may be used to construct instruments either in the reduced form equations or the instruments themselves. Unlike the standard linear models, the main disadvantage in the stochastic frontier setting is that a substantive assumption needs to be made regarding the correct specification of the reduced form in order to correctly predict the technical inefficiency component. In addition, the instruments, if they are available, often subject to potential pitfalls because they fail to meet the two required conditions adequately that the instruments are sufficiently correlated with the endogenous regressors, and they

1

are uncorrelated with the composed errors term. Thus, the potential difficulty of implementing these approaches is when there is no outside information available to construct the appropriate instruments.

To alleviate these problems, this paper considers an alternative approach to handle endogeneity in stochastic frontier models, which *does not require the availability of outside information to construct the instruments*. Consequently, we can construct a flexible joint distribution of the endogenous regressor and the composed error that can accommodate any degree of dependency between them. We then use this joint distribution to derive the likelihood function and maximize it to obtain the consistent estimates of the model parameters.

## 2. The Model and Methodology

Consider the following stochastic frontier model:

$$y_i = z_i' a + x_i' b + v_i - u_i, \quad i = 1, \mathrm{K}, n, \tag{1}$$

where $y_i$ is the output of firm $i$, $z_i$ is a $d \times 1$ vector of exogenous input, $x_i$ is a $p \times 1$ vector endogenous input, $a$ and $b$ are $d \times 1$ and $p \times 1$ vectors of unknown parameters, $v_i$ is a symmetric random error, $u_i$ is the one-sided random disturbance representing technical inefficiency. We assume that $z_i$ is uncorrelated with $v_i$ and $u_i$ but $x_i$ are allowed to be correlated with $v_i$ and possibly with $u_i$, and this generates the endogeneity problem. We also assume that $u_i$ and $v_i$ are independent and leave the form of $u_i$ unrestricted. The discussion that follows can be easily extended for the case where (exogenous) environmental variables are included in the distribution of $u_i$ (e.g., Battese and Coelli (1995)). Following standard practice, assume that $v_i : i.i.d. N(0, s_v^2)$ and $u_i : i.i.d. \left| N(0, s_u^2) \right|$. Then the density of $e_i = v_i - u_i = y_i - z_i' a - x_i' b$ is given by

$$g(e_i) = \int_0^{\Psi} f_v(e_i + u_i) f_u(u_i) \, du_i = \frac{2}{s} f \left( \frac{e_i}{s} \right) F \left( -\frac{l \, e_i}{s} \right) \tag{2}$$

where $s^2 = s_v^2 + s_u^2$, $l = s_u / s_v$, $f(.)$ and $F(.)$ are the probability density function and cumulative distribution function of a standard normal random variable, respectively.

Let $F(x_1, \mathrm{K}, x_p, e)$ and $f(x_1, \mathrm{K}, x_p, e)$ be the joint distribution and the joint density of $(x_1, \mathrm{K}, x_p)$ and $e_i$, respectively. In practice, $F(x_1, \mathrm{K}, x_p, e)$ and $f(x_1, \mathrm{K}, x_p, e)$ are typically unknown and hence need to be estimated. Following Park and Gupta (2012), we suggest a copula approach to construct and estimate this joint density. The copula essentially captures the dependence in the joint distribution of the endogenous regressors and the composed errors. For exposition purpose, suppose we have a joint distribution of $(x_1, \mathrm{K}, x_p, e)$ with joint density $f(x_1, \mathrm{K}, x_p, e)$, and let $f_j(x_j)$, $F_j(x_j)$,

2

for $j = 1, K, p$, $g(e)$ and $G(e)$ denote the marginal density and CDF of $x_j$ and $e$, respectively. Also let $C$ denotes the "copula function" defined for $(x_1, K, x_{p+1}) \hat{1} [0, 1]^{p+1}$ by

$$C(x_1, K, x_{p+1}) = P(F_1(x_1) \pounds x_1, K, F_p(x_p) \pounds x_p, G(e) \pounds x_{p+1}),$$

so that the copula function is itself a CDF. Moreover, since $F_j(x_j)$ and $G(.)$ are marginal distribution function, each component $U_j = F_j(x_j)$ and $U_e = G(e)$ has a uniform marginal distribution (see for example Li and Racine (2007, Theorem A.2)). Let $c(x_1, K, x_p)$ denotes the pdf associated with $C(x_1, K, x_p)$, then by Sklar's theorem (Sklar (1959)), we have

$$f(x_1, K, x_p, e) = c(F_1(x_1), K, F_p(x_p), G(e)) g(e) \mathop{\bigcirc}\limits_{j=1}^{p} f_j(x_j). \tag{3}$$

Thus, equation (3) shows that the copula function completely characterizes the dependence structure of $(x_1, K, x_p, e)$, and $c(x_1, K, x_p) = 1$ if and only if $(x_1, K, x_p, e)$ are independent of each other. For more rigorous treatment on Copula, see Nelsen (2006). To obtain the joint density in (3), we need to specify the copula function. One commonly used copula function is the Gaussian copula. Other copula functions such as Frank, Placket, Clayton, and Farlie-Gumbel-Morgenstern can also be used. The Gaussian copula is generally robust for most application and has many desirable properties (Danaher and Smith (2011)). Let $F_{S, p+1}$ denote a $(p+1)$-dimensional CDF with zero mean and correlation matrix $S$. Then the $(p+1)$-dimensional CDF with correlation matrix $S$ is given by $C(w, S) = F_{S, p+1}(F^{-1}(U_1), K, F^{-1}(U_p), F^{-1}(U_e))$, where $w = (U_1, K, U_p, U_e) = (F_1(x_1), K, F_p(x_p), G(e))$. The copula density is

$$c(w, S) = (\det(S))^{-1/2} \times$$
$$\exp\left[ -\frac{1}{2}(F^{-1}(U_1), K, F^{-1}(U_p), F^{-1}(U_e))'(S^{-1} - I_{p+1})(F^{-1}(U_1), K, F^{-1}(U_p), F^{-1}(U_e)) \right]. \tag{4}$$

The log-likelihood function corresponding to (5) is then

$$\ln L(q, S) = \mathop{\mathring{a}}\limits_{i=1}^{n} \left[ \ln c(F_1(x_{1i}), K, F_p(x_{pi}), G(e_i; q); S) + \mathop{\mathring{a}}\limits_{j=1}^{p} \ln f_j(x_{ji}) + \ln g(e_i; q) \right], \tag{5}$$

where $q = (a', b', l, s^2)'$ and the form of $c(.)$ is given in (4). Notice that the first term in the summation in (5) is derived from the copula density and this term reflects the dependence between the endogenous variables and the composed errors. In addition, since the marginal density $f_j(x_j)$ does not contain any parameters of interest, the second term in the summation in (5) can be dropped from the log-likelihood function. Finally, it is clear from (5) that if there are no endogeneity problem, (5) collapses to the log-likelihood function of the standard stochastic frontier models. By

3

maximizing the log-likelihood function in (5), consistent estimates of $(q, S)$ can be obtained, and this can be done in a two-step estimation procedure describe below.

_Step 1: Estimation of $F_j(x_j)$, $j = 1, K, p$; and $G(e, q)$_

Since we have observed sample of $x_{ji}$, $j = 1, K, p$; $i = 1, K, n$; in the first step, we can estimate $F_j(x_{ji})$ by

$$\hat{F}_{nj} = \frac{1}{n+1} \sum_{i=1}^{n} 1(x_{ji} \pounds x_{0j}), \quad j = 1, K, p, \tag{6}$$

where $1(.)$ is an indicator function. Note that in (6), we have used the rescaling factor $1/(n+1)$ rather than $1/n$ to avoid difficulties arising from the potential unboundedness of the $\ln c(F_1(x_{1i}), K, F_p(x_{pi}), G(e_i; q); S)$ as some of the $F_j(x_j)$ tend to one. To estimate $G(e_i; q)$, note that its density $g(e_i; q)$ is given in (2) and by definition, $G(e_i; q) = \int_{-\yen}^{e_i} g(s; q) ds$, thus $G(e_i; q)$ can be estimated using numerical integration, and let $\hat{G}(e_i; q)$ denotes the estimator of $G(e_i; q)$.

_Step 2: Maximize the log-likelihood function_

Maximize the log-likelihood function in (5) with $F_j(x_j)$ and $G(e_i; q)$ are replaced by their estimates $\hat{F}_j(x_j)$ and $\hat{G}(e_i; q)$, respectively, i.e.,

$$(\hat{q}, \hat{S}) = \underset{q \hat{I} Q, S}{\arg\max} \sum_{i=1}^{n} \left\{ \ln c(\hat{F}_1(x_{1i}), K, \hat{F}_p(x_{pi}), \hat{G}(e_i; q); S) + \ln g(e_i; q) \right\}. \tag{7}$$

_Predicting Technical Inefficiency:_

Once the parameters have been estimated, technical inefficiency $u_i$, can be predicted based on Jondrow et al. (1982):

$$\hat{u}_i = \hat{E}(u_i \mid e_i) = \frac{\hat{s}\hat{I}}{1 + \hat{I}^2} \left[ \frac{f(\hat{I}\hat{e}_i / \hat{s})}{1 - F(\hat{I}\hat{e}_i / \hat{s})} - \frac{\hat{I}\hat{e}_i}{\hat{s}} \right]$$

where $\hat{e}_i = y_i - z_i'\hat{a} - x_i'\hat{b}$ and $\hat{a}$, $\hat{b}$, $\hat{I}$ and $\hat{s}^2$ are the parameter estimates obtained from the Copula approach discussed above.

## 3. Monte Carlo Simulations

To examine the finite sample performance of the proposed Copula estimator we conduct some Monte Carlo experiments. We consider the following data generating process:

$$y_i = z_{1,i}a + x_i b + v_i - u_i,$$
$$x_i = z_{2,i}g + e_i,$$

where $u_i : i.i.d. \left| N(0, s_u^2) \right|$ and the random variables $z_{1,i}$ and $z_{i,2}$ are each generated independently as $c_{(2)}^2$. The vector of random errors $(v_i, e_i)'$ is generated by

$$\begin{pmatrix} v_i \\ e_i \end{pmatrix} : N\left(0, \begin{pmatrix} s_v^2 & rs_v s_e \\ rs_v s_e & s_e^2 \end{pmatrix}\right)$$

In our experiment, we fix $a = b = 0.5$, $s_e^2 = 1$, $s_v^2 = s_u^2 = 1$ and $g = 1$. We set the values of $r = \{0.0, 0.4, 0.8\}$ and consider two sample sizes: $n = (750, 1500)$. The simulations are replicated 1,000 times.[1] Note that since our DGP contains only one endogenous regressor, the Gaussian copula function has a simple form:

$$C(U_x, U_e) = N_{2,r}(F^{-1}(U_x), F^{-1}(U_e))$$

$$= \frac{1}{2p(1-r^2)^{1/2}} \int_{-\infty}^{F^{-1}(U_x)} \int_{-\infty}^{F^{-1}(U_e)} \exp\left\{\frac{-(s^2 + t^2 - 2rst)}{2(1-r^2)}\right\} ds dt,$$

where $N_{2,r}(.,.)$ denotes the standard bivariate normal distribution function with correlation $r$. The corresponding copula density is:

$$c(U_x, U_e) = \frac{1}{(1-r^2)^{1/2}} \exp\left\{\frac{-r^2(x^{*2} + e^{*2}) + 2rx^* e^*}{2(1-r^2)}\right\},$$

and the log-likelihood function associated with (5) can be expressed as:

$$\ln L(q) = -\frac{n}{2}\ln(1-r^2) - \sum_{i=1}^{n}\left\{\frac{r^2(x_i^{*2} + e_i^{*2}) - 2rx_i^* e_i^*}{2(1-r^2)}\right\} + \sum_{i=1}^{n} \ln g(y_i - z_i a - x_i b),$$

where $q = (a, b, r, l, s^2)'$, $x_i^* = F^{-1}(F(x_i))$, $e_i^* = F^{-1}(G(y_i - z_i a - x_i b))$ and the form of $g(.)$ is given in (2). In our simulation, the numerical integration to obtain $G(.)$ is performed using Gaussi-Kronrod quadrature with 50 points.

For comparison purpose, we also compute the standard MLE and GMM estimators of Tran and Tsionas (2013). Simulation results of the parameter estimates' MSE are displayed in Tables 1-3. Our simulations show that when there is no correlation (i.e., no endogeneity in the regressor), as expected, all estimators performed as well as the standard MLE for all ranges of the parameters considered. However, as expected, when there is correlation and as the correlation increases, the MLE deteriorates quickly and becomes severely biased, regardless of the sample sizes. Comparing the performance of the proposed Copula estimator to the GMM estimator of Tran and Tsionas (2013) shows that the proposed estimator perform quite well in term of MSE. Finally, as the sample size $n$ doubles, the estimated MSEs of the proposed Copula estimator reduces to about half of the original values regardless whether there is correlation or not; this is consistent with the fact that the Copula estimator is $\sqrt{n}$- consistent estimator.

## 4. Conclusion

---

[1] All computations are performed in Fortran 77 using extensively IMSL routines.

In this paper, we offer an alternative approach for estimating stochastic frontier models with endogenous regressors when there is no additional data available that is based on copula method to directly construct the joint density of the endogenous regressors and the composed errors to capture their dependency. We examine the finite sample behavior of the proposed approach via Monte Carlo simulations. The results from the simulations showed that the estimators are performed very well in finite samples in term of bias and MSE.

# References

Amsler, C., A. Prokhorov and P. Schmidt (2014), "Endogeneity in stochastic frontier models." *Working Paper*, Michigan State University.

Battese. G.E. and T.J. Coelli (1995), "A model for technical inefficiency effects in a stochastic frontier production function for panel data." *Empirical Economics*, 20, 325-332.

Danaher, P.J. and M.S. Smith (2011), "Modeling multivariate distributions using copulas: Applications in marketing." *Marketing Science*, 30 (1), 4-21.

Jondrow, J., C.A.K. Lovell, I.S. Materov and P. Schmidt (1982). "On the estimation of technical inefficiency in the stochastic frontier production function model." *Journal of Econometrics*, 19 (2/3), 233-238.

Kutlu, L. (2010), "Battese-Coelli estimator with endogenous regressors." *Economics Letters*, 109, 79-81.

Li, Q. and J. Racine, (2007). *Nonparametric Econometrics*. Princeton University Press, Princeton, NJ.

Nelsen, R.B. (2006), *An Introduction to Copula*, Vol. 139 of Springer Series in Statistics, Springer.

Park, S. and S. Gupta (2012), "Handling endogenous regressors by joint estimation using Copulas." *Marketing Science*, 31 (4), 567-586.

Sklar, A. (1959), "Fonctions de repartition a n dimensions et leurs marges." *Publications de l'Institut de Statistique de l'Uinversite de Paris,* 8:229-231.

Tran, K.C. and E.G. Tsionas (2013), "GMM estimation of stochastic frontier models with endogenous regressors." *Economics Letters*, 118, 233-236.

**Table 1: Estimated MSE:** $r = 0.0$, $s_u = 1$, $s_v = 1$

| | $n = 750$ | | | $n = 1500$ | | |
|---|---|---|---|---|---|---|
| | **MLE** | **GMM** | **Copula** | **MLE** | **GMM** | **Copula** |
| $a$ | .001577 | .001580 | .001585 | .000727 | .000729 | .000731 |
| $b$ | .000797 | .000797 | .000812 | .000404 | .000407 | .000409 |
| $g$ | - | .001099 | - | - | .000578 | - |
| $s_v$ | .001463 | .001458 | .001533 | .000735 | .000737 | .000736 |
| $s_u$ | .003002 | .003015 | .003241 | .001336 | .001339 | .001341 |
| $r$ | - | .003366 | .003561 | - | .001714 | .001717 |
| Total MSE[*] | .006839 | .006850 | .00717 | .003202 | .003212 | .003217 |

**[*]** Indicates the total MSE of the four parameters (as in the MLE case).


**Table 2: Estimated MSE:** $r = 0.4$, $s_u = 1$, $s_v = 1$

| | $n = 750$ | | | $n = 1500$ | | |
|---|---|---|---|---|---|---|
| | **MLE** | **GMM** | **Copula** | **MLE** | **GMM** | **Copula** |
| $a$ | .013767 | .001416 | .001543 | .007546 | .000728 | .000721 |
| $b$ | .033650 | .001579 | .001583 | .032991 | .000757 | .000763 |
| $g$ | - | .001149 | - | - | .000578 | - |
| $s_v$ | .003230 | .001638 | .001644 | .002206 | .000820 | .000815 |
| $s_u$ | .002927 | .002517 | .002510 | .001320 | .001193 | .001199 |
| $r$ | - | .002608 | .002600 | - | .001265 | .001273 |
| Total MSE[*] | .053574 | .007150 | .007155 | .044063 | .003497 | .003498 |

**[*]** See Table 1.


**Table 3: Estimated MSE:** $r = 0.8$, $s_u = 1$, $s_v = 1$

| | $n = 750$ | | | $n = 1500$ | | |
|---|---|---|---|---|---|---|
| | **MLE** | **GMM** | **Copula** | **MLE** | **GMM** | **Copula** |
| $a$ | .016041 | .001532 | .001542 | .008756 | .000708 | .000716 |
| $b$ | .133403 | .001296 | .001299 | .132104 | .000742 | .000747 |
| $g$ | - | .001037 | - | - | .000569 | - |
| $s_v$ | .026804 | .006471 | .006472 | .025901 | .001041 | .001044 |
| $s_u$ | .008483 | .001099 | .001095 | .001230 | .000692 | .000689 |
| $r$ | - | .000546 | .000543 | - | .000288 | .000286 |
| Total MSE[*] | 0.184731 | .010398 | .010408 | .167991 | .003183 | .003196 |

**[*]** See Table 1.