# EXACT DISTRIBUTIONS OF $R^2$ AND ADJUSTED $R^2$ IN A LINEAR REGRESSION MODEL WITH MULTIVARIATE $t$ ERROR TERMS

Kazuhiro Ohtani* and Hisashi Tanizaki*

In this paper we consider a linear regression model when error terms obey a multivariate $t$ distribution, and examine the effects of departure from normality of error terms on the exact distributions of the coefficient of determination (say, $R^2$) and adjusted $R^2$ (say, $\overline{R}^2$). We derive the exact formulas for the density function, distribution function and $m$-th moment, and perform numerical analysis based on the exact formulas. It is shown that the upward bias of $R^2$ gets serious and the standard error of $R^2$ gets large as the degrees of freedom of the multivariate $t$ error distribution (say, $\nu_0$) get small. The confidence intervals of $R^2$ and $\overline{R}^2$ are examined, and it is shown that when the values of $\nu_0$ and the parent coefficient of determination (say, $\Phi$) are small, the upper confidence limits are very large, relative to the value of $\Phi$.

*Key words and phrases*:  Adjusted $R^2$, exact distribution, interval estimation, multivariate $t$ error terms, $R^2$.

## 1.  Introduction

To measure goodness of fit of an estimated linear regression model, the coefficient of determination (say, $R^2$) and the adjusted coefficient of determination (say, $\overline{R}^2$) have traditionally been used (see Section 2 for $R^2$ and $\overline{R}^2$). Thus, there are many studies on the small sample properties of $R^2$ and $\overline{R}^2$. For example, Barten (1962) suggests a modified version of $R^2$ to reduce its bias, and Press and Zellner (1978) discuss the reason why the study of $R^2$ is important in the case of fixed regressors and they perform Bayesian analysis of $R^2$. Also, Cramer (1987) derives the exact formulas for the first two moments of $R^2$ and $\overline{R}^2$, and shows that $R^2$ is seriously biased upward in small samples while $\overline{R}^2$ is more unreliable than $R^2$ in terms of standard deviation.

Although it is assumed that the model is correctly specified in the above studies, Carrodus and Giles (1992) examine the small sample properties of $R^2$ when the independence of error terms is mistakenly assumed. Also, using asymmetric linear loss functions, Ohtani (1994) examines the risk performances of $R^2$ and $\overline{R}^2$ when the relevant regressors are omitted in the specified model and when irrelevant regressors are included in the specified model. Ohtani and Hasegawa (1993) examine the bias and mean squared error (MSE) performances when the proxy variables are used instead of unobservable regressors and the error terms obey a multivariate $t$ distribution.

Although there are many studies on the small sample properties of $R^2$ and $\overline{R}^2$, the studies on the exact distribution of $R^2$ and $\overline{R}^2$ per se are few. Although Ohtani (1994) derives the exact distribution and density functions of $R^2$ and $\overline{R}^2$, he assumes that the error terms obey a normal distribution. As is discussed in Fama (1965) and Blattberg and Gonedes (1974), there exist many economic data that may be generated by distributions with fatter tails than a normal distribution. One example of such distributions is a multivariate $t$ distribution. To examine the effects of departure from normality of error terms on the sampling performances of estimators and test statistics, the multivariate $t$ distribution has often been used. Some examples are Zellner (1976), Ullah and Zinde-Walsh (1984), Giles (1991), and Namba and Ohtani (2002). Although Srivastava and Ullah (1995) examined the sampling properties of $R^2$ and $\overline{R}^2$ under a general non-normal error distribution, their analysis is based on the large sample asymptotic expansions.

In this paper we consider a linear regression model when error terms obey a multivariate $t$ distribution, and examine the effects of departure from normality of error terms on the exact distributions of $R^2$ and $\overline{R}^2$. In Section 2 the model and estimators are presented, and in Section 3 the exact formulas for the density function, distribution function and $m$-th moment are derived. In Section 4 we evaluate means, standard errors, density functions, and confidence intervals of $R^2$ and $\overline{R}^2$ numerically. The numerical results show that the upward bias of $R^2$ gets serious and the standard error of $R^2$ gets large as the degrees of freedom of the multivariate $t$ error distribution (say, $\nu_0$) get small. It is also shown that when the values of $\nu_0$ and the parent coefficient of determination (say, $\Phi$, which is defined in Section 4) are small, the upper confidence limits of $R^2$ and $\overline{R}^2$ are vary large. Finally, the 95% confidence intervals of $R^2$ for $\Phi = 0.5$ are shown.

## 2. Model and estimators

We consider the following linear regression model:

$$(2.1) \qquad\qquad y = \ell\beta_0 + X\beta + u,$$

where $y$ is an $n \times 1$ vector of observations of the dependent variable, $\ell$ is an $n \times 1$ vector consisting of ones, $X$ is an $n \times (k-1)$ matrix of non-stochastic regressors, $\beta_0$ is an intercept, $\beta$ is a $(k-1) \times 1$ vector of regression coefficients, and $u$ is an $n \times 1$ vector of error terms.

As to the error terms, we assume that $u$ obeys a multivariate $t$ distribution with location parameter 0, scale parameter $\sigma^2$, and degrees of freedom parameter $\nu_0$. Then, as is shown in Zellner (1976), the density function of $u$ is written as:

$$(2.2) \qquad\qquad p(u) = \int_0^\infty p_N(u \,|\, \tau) \, p_{IG}(\tau) \, \mathrm{d}\tau,$$

where

$$(2.3) \qquad p_N(u|\tau) = \frac{1}{(2\pi)^{n/2}\,\tau^n} \exp\left(-\frac{u'u}{2\tau^2}\right),$$

$$(2.4) \qquad p_{IG}(\tau) = \frac{2}{\Gamma(\nu_0/2)} \left(\frac{\nu_0\sigma^2}{2}\right)^{\nu_0/2} \tau^{-(\nu_0+1)} \exp\left(-\frac{\nu_0\sigma^2}{2\tau^2}\right).$$

We assume that $\nu_0 > 2$ so that the first two moments of $u$ may exist. Then, we have $E[u] = 0$ and $E[uu'] = \sigma_u^2\,I_n = [\nu_0/(\nu_0 - 2)]\,\sigma^2\,I_n$.

We assume without loss of generality that all the regressors are measured as deviations from their sample means (i.e., $X'\ell = 0$). Then, the ordinary least squares (OLS) estimators of $\beta_0$ and $\beta$ are:

$$(2.5) \qquad\qquad\qquad \widehat{\beta}_0 = \ell'y/n = \overline{y},$$

$$(2.6) \qquad\qquad\qquad \widehat{\beta} = S^{-1}X'y,$$

where $S = X'X$. The associated residual vector is:

$$(2.7) \qquad\qquad\qquad e = y - (\ell\overline{y} + X\widehat{\beta}).$$

Since $y'y - n\overline{y}^2 = \widehat{\beta}'S\widehat{\beta} + e'e$, the sample coefficient of determination is written as:

$$(2.8) \qquad\qquad R^2 = 1 - \frac{e'e}{y'y - n\overline{y}^2} = \frac{\widehat{\beta}'S\widehat{\beta}}{\widehat{\beta}'S\widehat{\beta} + e'e}.$$

Also, the adjusted coefficient of determination is:

$$(2.9) \qquad\qquad \overline{R}^2 = 1 - \frac{n-1}{n-k}\left(1 - R^2\right).$$

If we define a formally general estimator as:

$$(2.10) \qquad\qquad R_h^2 = hR^2 + 1 - h,$$

where $h \geq 1$, then $R_h^2$ reduces to $R^2$ when $h = 1$, and to $\overline{R}^2$ when $h = (n-1)/(n-k)$. Since $0 \leq R^2 \leq 1$, we see that $1 - h \leq R_h^2 \leq 1$.

## 3. Exact density and distribution functions

If we assume temporarily that $\tau$ is fixed, then the error terms obey a normal distribution with $E[u] = 0$ and $E[uu'] = \tau^2\,I_n$. As is shown in Ohtani (1994), the density function of $R_h^2$, given $\tau$, is:

$$(3.1) \qquad p(c\,|\,\tau) = \sum_{i=0}^{\infty} \frac{w_i(\lambda)}{B((k-1)/2 + i, (n-k)/2)} h^{-(n-1)/2-i+1}$$

$$\times (c + h - 1)^{(k-1)/2+i-1}(1-c)^{(n-k)/2-1},$$

where $w_i(\lambda) = [(\lambda/2)^i/i!]\exp(-\lambda/2)$ and $\lambda = \beta'S\beta/\tau^2$.

Using (2.4) and (3.1), the density function of $R_h^2$ can be obtained as follows:

$$
(3.2) \qquad p(c) = \int_0^\infty p(c\,|\,\tau)\,p_{IG}(\tau)\,\mathrm{d}\tau
$$

$$
= \sum_{i=0}^\infty \frac{1}{B((k-1)/2+i,(n-k)/2)\,i!}h^{-(n-1)/2-i+1}
$$

$$
\times (c+h-1)^{(k-1)/2+i-1}(1-c)^{(n-k)/2-1}
$$

$$
\times \frac{2}{\Gamma(\nu_0/2)}\left(\frac{\nu_0\sigma^2}{2}\right)^{\nu_0/2}\left(\frac{\beta'S\beta}{2}\right)^i
$$

$$
\times \int_0^\infty \tau^{-(\nu_0+1)-2i}\exp\left(-\frac{\beta'S\beta+\nu_0\sigma^2}{2\tau^2}\right)\,\mathrm{d}\tau.
$$

Making use of the change of variable, $t = (\beta'S\beta + \nu_0\sigma^2)/(2\tau^2)$, and performing some manipulations, we obtain the following distribution:

$$
(3.3) \qquad p(c) = \sum_{i=0}^\infty \frac{\theta^i\,\nu_0^{\nu_0/2}\,\Gamma(\nu_0/2+i)}{B((k-1)/2+i,(n-k)/2)\,i!\,\Gamma(\nu_0/2)(\nu_0+\theta)^{\nu_0/2+i}}
$$

$$
\times h^{-(n-3)/2-i}(c+h-1)^{(k-1)/2+i-1}(1-c)^{(n-k)/2-1},
$$

where $\theta = \beta'S\beta/\sigma^2$, and $B(\cdot,\cdot)$ is the beta function.

The distribution function of $R_h^2$ is:

$$
(3.4)
$$

$$
F(c_0) = \int_{1-h}^{c_0} p(c)\,\mathrm{d}c
$$

$$
= \sum_{i=0}^\infty \frac{\theta^i\,\nu_0^{\nu_0/2}\,\Gamma(\nu_0/2+i)}{B((k-1)/2+i,(n-k)/2)\,i!\,\Gamma(\nu_0/2)(\nu_0+\theta)^{\nu_0/2+i}}
$$

$$
\times h^{-(n-3)/2-i}\int_{1-h}^{c_0}(c+h-1)^{(k-1)/2+i-1}(1-c)^{(n-k)/2-1}\,\mathrm{d}c.
$$

Making use of the change of variable, $t = (c+h-1)/h$, and performing some manipulations, (3.4) reduces to:

$$
(3.5) \qquad F(c_0) = \sum_{i=0}^\infty \frac{\theta^i\,\nu_0^{\nu_0/2}\,\Gamma(\nu_0/2+i)}{i!\,\Gamma(\nu_0/2)(\nu_0+\theta)^{\nu_0/2+i}}\,I_{c_0^*}((k-1)/2+i,(n-k)/2).
$$

where $c_0^* = (c_0+h-1)/h$, and $I_a(\cdot,\cdot)$ is the incomplete beta function ratio. When $\beta = 0$ (i.e., $\theta = 0$), we see that the distribution function reduces to:

$$
(3.6) \qquad\qquad F(c_0) = I_{c_0^*}((k-1)/2,(n-k)/2).
$$

Putting $\lambda_1 = \lambda_2 = 0$ in eq. (18) in Ohtani (1994) which is the distribution function when the error terms obey a normal distribution, and comparing with

(3.6), we see that when $\beta = 0$, the distribution function is robust to the change of the error distribution from a normal distribution to a multivariate $t$ distribution. However, when $\beta \neq 0$, this robustness does not hold.

Also, the formula for the $m$-th moment of $R_h^2$ is:

(3.7)

$$
\mathrm{E}\left[(R_h^2)^m\right] = \int_{1-h}^1 c^m \, p(c) \, \mathrm{d}c
$$

$$
= \sum_{i=0}^\infty \frac{\theta^i \, \nu_0^{\nu_0/2} \, \Gamma(\nu_0/2 + i)}{B((k-1)/2 + i, (n-k)/2) \, i! \, \Gamma(\nu_0/2)(\nu_0 + \theta)^{\nu_0/2+i}}
$$

$$
\times \, h^{-(n-3)/2-i} \int_{1-h}^1 c^m \, (c+h-1)^{(k-1)/2+i-1} \, (1-c)^{(n-k)/2-1} \, \mathrm{d}c.
$$

Again, making use of the change of variable, $t = (c+h-1)/h$, the integral in (3.7) reduces to:

$$
(3.8) \quad \int_0^1 [th + (1-h)]^m \, (th)^{(k-1)/2+i-1} \, [(1-t)h]^{(n-k)/2-1} \, h \, \mathrm{d}t
$$

$$
= \sum_{r=0}^m {}_mC_r \, h^{(n-3)/2+r+i} \, (1-h)^{m-r} \int_0^1 t^{(k-1)/2+r+i-1} \, (1-t)^{(n-k)/2-1} \, \mathrm{d}t
$$

$$
= \sum_{r=0}^m {}_mC_r \, h^{(n-3)/2+r+i} \, (1-h)^{m-r} \, B((k-1)/2 + r + i, (n-k)/2).
$$

Thus, using the formula, $B(a, b) = \Gamma(a)\Gamma(b)/\Gamma(a+b)$, we finally obtain the expectation of $(R_h^2)^m$:

$$
(3.9) \qquad \mathrm{E}\left[(R_h^2)^m\right] = \sum_{i=0}^\infty \frac{\Gamma((n-1)/2 + i) \, \Gamma(\nu_0/2 + i) \, \theta^i \, \nu_0^{\nu_0/2}}{\Gamma((k-1)/2 + i)\Gamma(\nu_0/2) \, i! \, (\nu_0 + \theta)^{\nu_0/2+i}}
$$

$$
\times \sum_{r=0}^m {}_mC_r h^r \, (1-h)^{m-r} \frac{\Gamma((k-1)/2 + r + i)}{\Gamma((n-1)/2 + r + i)}.
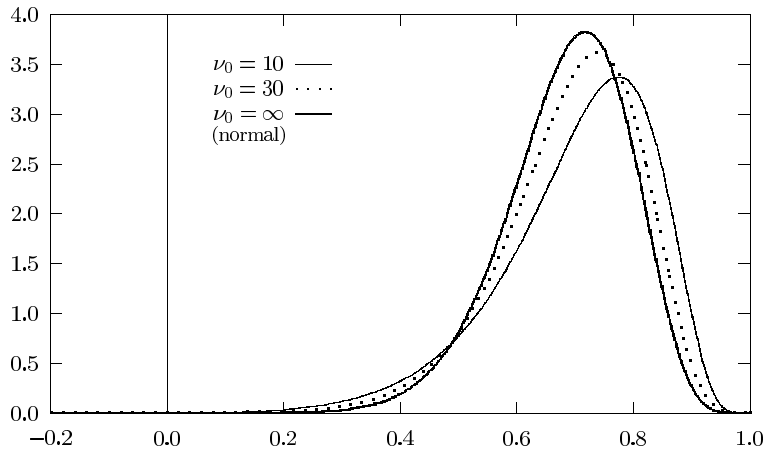$$

## 4. Numerical analysis

In this section we perform numerical analysis based on the exact formulas given in (3.3), (3.5) and (3.9). We define $\Phi$ as follows:

$$
(4.1) \qquad\qquad \Phi = \frac{\beta' S \beta}{\beta' S \beta + n\sigma_u^2} = \frac{\theta}{\theta + n\nu_0/(\nu_0 - 2)},
$$

which is called the parent coefficient of determination (see Press and Zellner (1978), Cramer (1987) and Ohtani and Hasegawa (1993)). Note that the relationship between $\Phi$ and $R^2$ is given by $\mathrm{plim}_{n\to\infty} \Phi = \mathrm{plim}_{n\to\infty} R^2$. In the numerical evaluations, we first decided the value of $\Phi$, and then calculated the value of $\theta$

Table 1.  Mean, standard error and 95% confidence interval of $R^2$ for $k = 5$ and $n = 20$.

| $\nu_0$ | $\Phi$ | Mean | S.E. | $c_L$ | $c_U$ |
|---|---|---|---|---|---|
| 5 | 0.8 | 0.8633 | 0.1028 | 0.5984 | 0.9731 |
|  | 0.6 | 0.7307 | 0.1525 | 0.3832 | 0.9298 |
|  | 0.4 | 0.5881 | 0.1774 | 0.2464 | 0.8688 |
|  | 0.2 | 0.4226 | 0.1762 | 0.1485 | 0.7684 |
| 10 | 0.8 | 0.8522 | 0.0823 | 0.6535 | 0.9563 |
|  | 0.6 | 0.7043 | 0.1334 | 0.4175 | 0.9014 |
|  | 0.4 | 0.5513 | 0.1613 | 0.2536 | 0.8297 |
|  | 0.2 | 0.3885 | 0.1640 | 0.1420 | 0.7307 |
| 30 | 0.8 | 0.8476 | 0.0650 | 0.7001 | 0.9439 |
|  | 0.6 | 0.6914 | 0.1163 | 0.4526 | 0.8802 |
|  | 0.4 | 0.5323 | 0.1493 | 0.2645 | 0.8051 |
|  | 0.2 | 0.3717 | 0.1574 | 0.1394 | 0.7114 |
| 100 | 0.8 | 0.8464 | 0.0583 | 0.7183 | 0.9387 |
|  | 0.6 | 0.6875 | 0.1094 | 0.4675 | 0.8722 |
|  | 0.4 | 0.5265 | 0.1448 | 0.2696 | 0.7969 |
|  | 0.2 | 0.3666 | 0.1554 | 0.1387 | 0.7056 |
| $\infty$ | 0.8 | 0.8459 | 0.0553 | 0.7206 | 0.9351 |
| (normal) | 0.6 | 0.6860 | 0.1063 | 0.4514 | 0.8633 |
|  | 0.4 | 0.5241 | 0.1429 | 0.2283 | 0.7796 |
|  | 0.2 | 0.3645 | 0.1545 | 0.0866 | 0.6732 |



Figure 1.  Density functions of $R^2$ for $k = 5$, $n = 20$, and $\Phi = 0.6$.

through $\theta = n\nu_0\Phi/[(\nu_0 - 2)(1 - \Phi)]$.  The parameter values used in the numerical evaluations were $k = 3, 4, 5, 6, 7, 8, n = 10, 20, 30, 40, \nu_0 = 3, 5, 10, 30, 100, \infty$ (normal), and various values of $\Phi$.  The numerical evaluations were executed on a personal computer, using the FORTRAN code.  The infinite series in the exact formulas converged rapidly with the convergence tolerance of $10^{-12}$.

Tables 1 and 2 show the mean, standard error (denoted as 'S.E.') and 95% confidence interval of $R^2$ and $\overline{R}^2$ when $k = 5$ and $n = 20$, where '$c_L$' and '$c_U$'

Table 2. Mean, standard error and 95% confidence interval of $\overline{R}^2$ for $k = 5$ and $n = 20$.

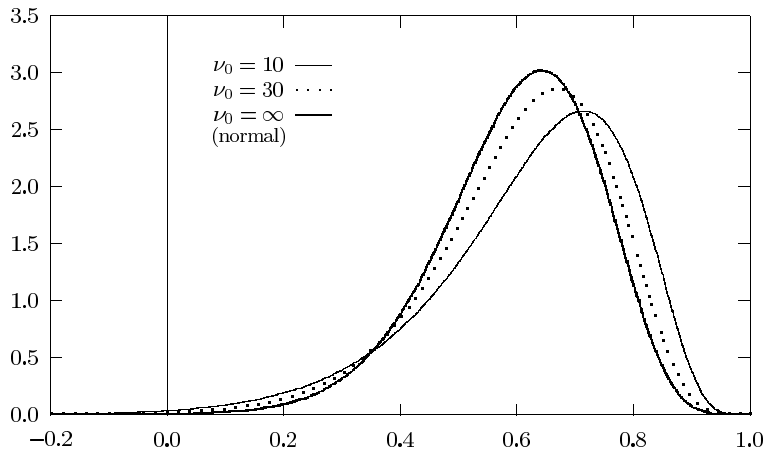| $\nu_0$ | $\Phi$ | Mean | S.E. | $c_L$ | $c_U$ |
|---|---|---|---|---|---|
| 5 | 0.8 | 0.8269 | 0.1302 | 0.4913 | 0.9659 |
| | 0.6 | 0.6589 | 0.1932 | 0.2187 | 0.9110 |
| | 0.4 | 0.4783 | 0.2247 | 0.0454 | 0.8338 |
| | 0.2 | 0.2687 | 0.2232 | $-0.0785$ | 0.7067 |
| 10 | 0.8 | 0.8128 | 0.1042 | 0.5611 | 0.9447 |
| | 0.6 | 0.6254 | 0.1689 | 0.2621 | 0.8751 |
| | 0.4 | 0.4316 | 0.2043 | 0.0546 | 0.7843 |
| | 0.2 | 0.2254 | 0.2078 | $-0.0868$ | 0.6589 |
| 30 | 0.8 | 0.8070 | 0.0823 | 0.6202 | 0.9290 |
| | 0.6 | 0.6091 | 0.1473 | 0.3066 | 0.8482 |
| | 0.4 | 0.4076 | 0.1891 | 0.0684 | 0.7532 |
| | 0.2 | 0.2041 | 0.1994 | $-0.0901$ | 0.6345 |
| 100 | 0.8 | 0.8055 | 0.0739 | 0.6432 | 0.9224 |
| | 0.6 | 0.6042 | 0.1386 | 0.3255 | 0.8382 |
| | 0.4 | 0.4002 | 0.1835 | 0.0748 | 0.7427 |
| | 0.2 | 0.1977 | 0.1968 | $-0.0910$ | 0.6270 |
| $\infty$ | 0.8 | 0.8049 | 0.0701 | 0.6461 | 0.9177 |
| (normal) | 0.6 | 0.6022 | 0.1346 | 0.3052 | 0.8268 |
| | 0.4 | 0.3972 | 0.1810 | 0.0226 | 0.7208 |
| | 0.2 | 0.1951 | 0.1957 | $-0.1570$ | 0.5860 |



Figure 2. Density functions of $\overline{R}^2$ for $k = 5$, $n = 20$, and $\Phi = 0.6$.

denote the confidence limits such that $P(R^2 < c_L) = P(\overline{R}^2 < c_L) = 0.025$ and $P(R^2 > c_U) = P(\overline{R}^2 > c_U) = 0.025$, where $P(A)$ is the probability of an event $A$. Figures 1 and 2 show the density functions of $R^2$ and $\overline{R}^2$ when $k = 5$, $n = 20$ and $\Phi = 0.6$.

We see from Tables 1 and 2 that $R^2$ is seriously biased upward in small samples, and $\overline{R}^2$ is more unreliable than $R^2$ in terms of standard error. In particular, the upward bias of $R^2$ gets serious and the standard error of $R^2$ gets

Table 3.　95% confidence interval when $\Phi = 0.5$.

| $k$ | $n$ | $\nu_0 = \infty$ | | $\nu_0 = 3$ | |
|---|---|---|---|---|---|
| | | $c_L$ | $c_U$ | $c_L$ | $c_U$ |
| 3 | 10 | 0.2770 | 0.8678 | 0.2933 | 0.9658 |
| | 20 | 0.3187 | 0.7558 | 0.2361 | 0.9374 |
| | 30 | 0.3445 | 0.7056 | 0.2170 | 0.9270 |
| | 40 | 0.3617 | 0.6760 | 0.2073 | 0.9295 |
| 4 | 10 | 0.3470 | 0.9008 | 0.3676 | 0.9741 |
| | 20 | 0.3489 | 0.7762 | 0.2743 | 0.9421 |
| | 30 | 0.3638 | 0.7199 | 0.2432 | 0.9302 |
| | 40 | 0.3758 | 0.6871 | 0.2273 | 0.9318 |
| 5 | 10 | 0.4205 | 0.9310 | 0.4443 | 0.9819 |
| | 20 | 0.3794 | 0.7963 | 0.3125 | 0.9468 |
| | 30 | 0.3831 | 0.7342 | 0.2694 | 0.9334 |
| | 40 | 0.3900 | 0.6981 | 0.2472 | 0.9342 |
| 6 | 10 | 0.4983 | 0.9575 | 0.5241 | 0.9889 |
| | 20 | 0.4104 | 0.8161 | 0.3508 | 0.9515 |
| | 30 | 0.4026 | 0.7483 | 0.2954 | 0.9366 |
| | 40 | 0.4041 | 0.7090 | 0.2671 | 0.9365 |
| 7 | 10 | 0.5822 | 0.9791 | 0.6083 | 0.9948 |
| | 20 | 0.4418 | 0.8354 | 0.3893 | 0.9562 |
| | 30 | 0.4222 | 0.7624 | 0.3215 | 0.9398 |
| | 40 | 0.4184 | 0.7199 | 0.2869 | 0.9388 |
| 8 | 10 | 0.6755 | 0.9940 | 0.6996 | 0.9987 |
| | 20 | 0.4737 | 0.8544 | 0.4279 | 0.9608 |
| | 30 | 0.4419 | 0.7763 | 0.3476 | 0.9430 |
| | 40 | 0.4327 | 0.7308 | 0.3067 | 0.9412 |

large as the degrees of freedom of the multivariate $t$ error distribution get small. The phenomena are also seen from Figure 1. This indicates that as the tails of the error distribution get fatter, $R^2$ becomes more unreliable. Also, we see from Figures 1 and 2 that the density function of $\overline{R}^2$ is flatter than that of $R^2$ though the modes of the density functions of $\overline{R}^2$ are smaller than those of $R^2$.

We see from Tables 1 and 2 that the confidence intervals of $R^2$ and $\overline{R}^2$ are considerably wide, and the confidence intervals get wide as the degrees of freedom of the multivariate $t$ error distribution get small. This phenomenon is also expected from Figures 1 and 2. In particular, when the values of $\nu_0$ and $\Phi$ are small, the upper confidence limits of $R^2$ and $\overline{R}^2$ are vary large. For example, when $\nu_0 = 5$ and $\Phi = 0.2$, the upper confidence limit of $R^2$ is $c_U = 0.7684$, and that of $\overline{R}^2$ is $c_U = 0.7067$. This indicates that even when the estimated values of $R^2$ and $\overline{R}^2$ are more than 0.7, the parent coefficient of determination may be just 0.2. We see from Table 2 that when the value of $\Phi$ is small, the lower confidence limits of $\overline{R}^2$ can be negative though the absolute value of $c_L$ becomes small. This phenomenon is caused by the shift to the right of the density function when $\nu_0$ decreases, as is shown in Figure 2.

Finally, we show 95% confidence intervals for $\Phi = 0.5$ and for some values of $k$ and $n$ in Table 3. Although there is no definite reason why $\Phi = 0.5$ is selected, we can confirm at the confidence coefficient 0.95 that the parent coefficient of

determination is at least more than half if the value of $R^2$ exceeds the upper limit given in Table 3. Since $\nu_0 = \infty$ and $\nu_0 = 3$ are two extreme values, we can confirm at least $\Phi = 0.5$ if the value of $R^2$ is larger than the upper limit for $\nu_0 = 3$ even if the true value of $\nu_0$ is larger than 3, and we may doubt $\Phi = 0.5$ if the value of $R^2$ is less than the lower limit for $\nu_0 = \infty$.

## Acknowledgements

## REFERENCES

Barten, A. P. (1962). Note in the unbiased estimation of the squared multiple correlation coefficient, *Satistica Neerlandica*, **16**, 151–163.

Blattberg, R. C. and Gonedes, N. J. (1974). A comparison of the stable and Student distributions as statistical models for stock prices, *Journal of Business*, **47**, 244–280.

Carrodus, M. L. and Giles, D. E. A. (1992). The exact distribution of $R^2$ when regression disturbances are autocorrelated, *Economics Letters*, **38**, 375–380.

Cramer, J. S. (1987). Mean and variance of $R^2$ in small and moderate samples, *Journal of Econometrics*, **35**, 253–266.

Fama, E. F. (1965). The behaviour of stock market prices, *Journal of Business*, **38**, 34–105.

Giles, J. A. (1991). Pre-testing for linear restrictions in a regression model with spherically symmetric disturbances, *Journal of Econometrics*, **50**, 377–398.

Namba, A. and Ohtani, K. (2002). MSE performance of the double $k$-class estimator of each individual regression coefficient under multivariate $t$-errors, *Handbook of Applied Econometrics and Statistical Inference*, (eds. Ullah, A., Wan, A. T. K. and Chaturvedi, A.), 305–326.

Ohtani, K. (1994). The density functions of $R^2$ and $\overline{R}^2$, and their risk performance under asymmetric loss in misspecified linear regression models, *Economic Modelling*, **11**, 463–471.

Ohtani, K. and Hasegawa, H. (1993). On small sample properties of $R^2$ in a linear regression model with multivariate $t$ errors and proxy variables, *Econometric Theory*, **9**, 504–515.

Press, S. J. and Zellner, A. (1978). Posterior distribution for the multiple correlation coefficient with fixed regressors, *Journal of Econometrics*, **8**, 307–321.

Srivastava, A. K. and Ullah, A. (1995). The coefficient of determination and its adjusted version in linear regression models, *Econometric Reviews*, **14**, 229–240.

Ullah, A. and Zinde-Walsh, V. (1984). On the robustness of LM, LR, and Wald tests in regression model, *Econometrica*, **52**, 1055–1066.

Zellner, A. (1976). Bayesian and non-Bayesian analysis of the regression model with multivariate Student-t error terms, *Journal of the American Statistical Association*, **71**, 400–405.