

# Missing Value Imputation for Non-Normal Data

Ravi Khattree

Distinguished University Professor of Applied Statistics  
Co-Director, Center for Data Sciences and Big Data Analytics  
Participating Member, Center for Biomedical Research  
Oakland University

(Joint work with Zhixin Lun and Manoj Bahuguna)  
email: [khattree@oakland.edu](mailto:khattree@oakland.edu)

Marquette U., Oct. 25, 2024

# Overview

- Introduction
- Multivariate Copula-Transformation
- Copula Transformed Multiple Regression
- Data Sets
- Transformation Method, Analysis and Evaluation
- Imputation of Missing Values via Copula Transformation-MCAR
- Concluding Remarks

# Introduction

Basic Result (Casella and Berger, Theorem 2.1.10, p. 54):

For a continuous random variable  $W$ , the cumulative distribution function  $U = F_W(\cdot)$  is uniformly distributed on  $[0,1]$ .

Thus, any continuous random variable  $W$  can be transformed to a uniform random variable and conversely, any uniform random variable can be transformed to a random variable  $V$  with any desired continuous distribution.

# Introduction

This fact, thus, provides an approach for the transformation of data to any desirable distribution.

In particular, this can be used as an all purpose approach to introduce normality so that analyses, such as linear model theory based modeling, can be performed under the usual normality assumptions.

NORTA (NORMAL to ANYTHING) algorithms use this fact backwards to simulate random numbers from any desired distribution after generating standard normal variates.

# Introduction

However, one fact that should not be lost in this approach is that the information about the location and scale is essentially lost, in much the same way as the information about the shape and the skewness of the original distribution.

While symmetrized (specifically normally distributed) data are desirable for data analysis, the approach will therefore, prevent us from making any inference about the location in most situations.

A common approach to deal with skewed data is through the use of transformation (to normality). The logarithmic, square-root, arcsine and more generally, Box-Cox transformations have been the common tools to artificially induce, among other desirable features, symmetry for the asymmetric data and have been used extensively in a variety of statistical problems.

# Introduction

In the multivariate context, when the interest is in studying the dependence structures and possibly prediction, a generalization of the approach described above can be extremely useful

and  
the purpose of this talk is to introduce the usefulness of our suggested approach in various multivariate situations.

The objective here is to come up with an approach to transform the data where classical techniques of multivariate analyses [**Specifically, here for missing data imputation**] can be readily adopted for the transformed data.

Yet, the method should be such, so that inference and especially the predictions for the transformed data can be brought back to the original context.

# Multivariate Copula-Transformation

We will first define the concept of copula.

## Definition

A function  $C$  from a  $d$ -dimensional rectangle  $[0, 1]^d$  to  $[0, 1]$  is called a copula if there is a random vector  $\mathbf{U} = (U_1, U_2, \dots, U_d)'$ , such that for  $i = 1, \dots, d$ ,  $U_i \sim U(0, 1)$ , the uniform distribution on interval  $[0, 1]$  and  $C(u_1, u_2, \dots, u_d) = P[U_1 \leq u_1, U_2 \leq u_2, \dots, U_d \leq u_d]$  where  $U_1, U_2, \dots, U_d \in [0, 1]$ .

Thus,  $C(\cdot)$  is essentially a  $d$ -dimensional multivariate cumulative distribution function of  $d$  random variables, each distributed uniformly in the interval  $[0, 1]$ .

The dependence structure is not stated in the definition and cannot be, in general, specified. It depends on the nature and the joint behavior of the particular set of the random variables.

# Multivariate Copula-Transformation

More light on this issue is shed by Sklar's Theorem. Assume our random variables to be all continuous valued.



# Multivariate Copula-Transformation

## Theorem

*(Sklar's Theorem)*

A function  $F : R^d \rightarrow [0, 1]$  is the distribution function of some continuous random vector  $\mathbf{X} = (X_1, X_2, \dots, X_d)'$  iff there is a copula  $C$  from  $[0, 1]^d$  to  $[0, 1]$  and  $d$  univariate distribution functions  $F_1, F_2, \dots, F_d$  such that

$$C(F_1(x_1), F_2(x_2), \dots, F_d(x_d)) = F(x_1, x_2, \dots, x_d) \quad (1)$$

for  $\mathbf{X} = (X_1, X_2, \dots, X_d)' \in R^d$ .

The functions  $F_i(\cdot)$  are clearly the (marginal) cumulative distribution functions of corresponding random variables  $X_i, i = 1, 2, \dots, d$ . Thus, copula expresses the dependence among  $X_1, X_2, \dots, X_d$  through their marginal cumulative distribution functions.

# Multivariate Copula-Transformation

It provides a way to express and obtain the joint cumulative distribution functions through an appropriate copula. Since  $F(\cdot)$  is continuous and hence admits an inverse function  $F^{-1}(\cdot)$ , it follows from above that

$$(X_1, X_2, \dots, X_d) \stackrel{dist}{=} F^{-1}C(F_1(x_1), F_2(x_2), \dots, F_d(x_d)), \quad (2)$$

where *dist* indicates the equivalence of probability distributions.

# Multivariate Copula-Transformation

From (??) it follows that since  $U_i = F_i(X_i)$  is uniformly distributed in interval  $[0,1]$ ,  $i = 1, 2, \dots, d$  and since  $u_i = F(x_i) \implies x_i = F^{-1}(u_i)$ , for  $\mathbf{u} = (u_1, u_2, \dots, u_d)'$ , we have,

$$\begin{aligned} C(\mathbf{u}) &= C(u_1, u_2, \dots, u_d) \\ &= F(F_1^{-1}(u_1), F_2^{-1}(u_2), \dots, F_d^{-1}(u_d)), \mathbf{u} \in [0, 1]^d. \end{aligned} \quad (3)$$

# Multivariate Copula-Transformation

Let us concentrate on (??) namely,

$$C(F_1(x_1), F_2(x_2), \dots, F_d(x_d)) = F(x_1, x_2, \dots, x_d)$$

Let  $F(\cdot)$  and  $G(\cdot)$  be two  $d$ -dimensional multivariate continuous CDFs, with corresponding marginal CDFs  $F_1(\cdot), F_2(\cdot), \dots, F_d(\cdot)$  and  $G_1(\cdot), G_2(\cdot), \dots, G_d(\cdot)$  respectively. Also assume that  $F(\cdot)$  and  $G(\cdot)$  both correspond to the same copula function  $C(\cdot)$ . Thus, with random vector  $\mathbf{X}$  having the CDF  $F(\cdot)$  and random vector  $\mathbf{Y}$  having that as  $G(\cdot)$ , we have

$$\begin{aligned} F(x_1, x_2, \dots, x_d) &= C(F_1(x_1), F_2(x_2), \dots, F_d(x_d)) & (4) \\ &= C(u_1, u_2, \dots, u_d) = C(G_1(y_1), G_2(y_2), \dots, G_d(y_d)) \\ &= G(y_1, y_2, \dots, y_d), \end{aligned}$$

for some  $y_1, y_2, \dots, y_d$ , so that  $G^{-1}(y_i) = u_i, i = 1, 2, \dots, d$  where  $G^{-1}(\cdot)$  is the inverse function of  $G(\cdot)$ .

# Multivariate Copula-Transformation

Note that the only assumption here is that the  $F(\cdot)$  and  $G(\cdot)$  share the common copula. By Sklar's Theorem, it also follows that given  $F(\cdot)$  (or  $G(\cdot)$ ), the copula is unique.

Thus, if  $G(\cdot)$  is desired to be a particular cumulative distribution function then it automatically determines the choice of  $C(\cdot)$ .

# Multivariate Copula-Transformation

The above calculation succinctly provides, an approach to transform the data on multivariate random vector  $\mathbf{X}$  having cumulative distribution function  $F(\cdot)$  to another random vector  $\mathbf{Y}$  having the cumulative distribution function  $G(\cdot)$ . More succinctly, it can be described, analogous to (??) as,

$$(Y_1, Y_2, \dots, Y_d) \stackrel{dist}{=} G^{-1}C(F_1(X_1), F_2(X_2), \dots, F_d(X_d)). \quad (5)$$

# Multivariate Copula-Transformation

This is pictorially depicted in Figure on next slide.

For our work, with an intention to enable us to do classical multivariate modeling, the function  $G(\cdot)$  will usually be a multivariate normal cumulative distribution function. Consequently, the choice of  $C(\cdot)$  must be a Gaussian copula.

# Multivariate Copula-Transformation

As a graphical representation for two given distribution functions say,  $\mathcal{F}(\cdot)$  and  $\mathcal{G}(\cdot)$  with common copula say  $\mathcal{C}(\cdot)$ , our transformation works as,

$$\mathcal{F}(\cdot) \longrightarrow \mathcal{C}(\cdot) \longrightarrow \mathcal{G}(\cdot)$$

Figure 1



# Multivariate Copula-Transformation

Accordingly, implicit assumption on the distribution function  $F(\cdot)$  of our raw data is that even though  $F(\cdot)$  itself may not be multivariate normal distribution function, its copula function is Gaussian.

Such an assumption is very reasonable.

Of course, for the data analysis, we must resort to the empirical version of  $F(\cdot)$ . computed from data.

# Multivariate Copula-Transformation

Thus, in essence, we make the assumption that the common copula is a Gaussian copula  $\Phi_{(\boldsymbol{\mu}, \boldsymbol{\Sigma})}(\cdot)$ ,

In principle, the choices of mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$  are arbitrary.

Since our interest is in doing the multivariate analyses of dependence, we will choose  $\boldsymbol{\Sigma}$  cautiously to retain the essential dependence features of data.

On the other hand, since the choice of  $\boldsymbol{\mu}$  is often unimportant in such situations, we will take it's value to be the zero vector.

# Multivariate Copula-Transformation

As a graphical representation for a given distribution of data, say  $\mathcal{D}$ , our transformation works as

$$\mathcal{D} \longrightarrow \mathcal{U} \longrightarrow \mathcal{N}$$

and in a reverse direction as

$$\mathcal{N} \longrightarrow \mathcal{U} \longrightarrow \mathcal{D},$$

where  $\mathcal{U}$  represents a multivariate distribution in which each marginal is uniform on  $[0,1]$  and  $\mathcal{N}$  represents the chosen multivariate normal distribution.

# Multivariate Copula-Transformation

As a result, the copula used is the Gaussian copula defined by,

$$C_{\Sigma}(u_1, \dots, u_d) = \Phi_{\Sigma}(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_d))$$

where  $\Phi(\cdot)$  is the distribution function of (univariate) standard normal random variable and  $\Phi_{\Sigma}(\cdot)$  is the  $d$ -variate standard normal cumulative distribution function with mean vector  $\mathbf{0}$  and covariance matrix  $\Sigma$ .

**Regression Modeling in Copula space is crucial for our missing data imputation approach. So we will illustrate that first.**

# Copula Transformed Multiple Regression: Data Sets

- i) Wicklin's Data (2013): Taken from Wicklin's book, where we have four random variables, jointly exhibiting dependence, but each with marginal distributions which are functionally very different.

Specifically, we have the response variable  $y$  distributed as standard lognormal ( $\mu = 0, \sigma = 1$ ) and explanatory variables  $x_1, x_2$  and  $x_3$  respectively, distributed as standard normal, uniform on  $[0,1]$  and standard exponential ( $\lambda = 1$ ).

Clearly, considerable skewness is present in  $y$  and  $x_3$ . Also, the conditional distribution of  $y$  given  $x_1, x_2$  and  $x_3$  is clearly not normal. A total of 100 observations are available.

# Copula Transformed Multiple Regression: Data Sets

Other Data sets (not discussed):

- ii) Financial Indexes Data
- iii) SENIC Data
- iv) Prostate Cancer Data
- v) Real Estate Sales Data
- vi) Used Car Data
- vii) University Admissions Data

# Copula Transformed Multiple Regression: Analyses

For all the above data set we will fit the multiple linear regression model regressing  $y$  on  $k$  explanatory variables  $x_1, x_2, \dots, x_k$ . Clearly, the value of  $k$  is different for various data sets. No cross product or higher degree polynomial terms are assumed.

The objective is to compare the regression models fitted on the original data with those obtained by fitting the equivalent model on the corresponding Gaussian-copula transformed data.

# Copula Transformed Multiple Regression: Analyses

Specifically, since the number of observations and the functional forms of the models will be the same in the two situations, the coefficient of determination  $R^2$  values can be compared, along with the statistical significance of the models.

However, from a practical point of view, quality of prediction is also important and thus, we will also compare the prediction errors as well as the prediction intervals.

**[After all we are going to predict the missing values so quality of prediction better be superior.]**

For a fair comparison, these two will be obtained for the original response variable.



# Copula Transformed Multiple Regression: Analyses

Here we naturally replace  $C = C_{\mathbf{R}}$  (Gaussian Copula),  $G = \Phi_{\mathbf{R}}$  (as we are using Gaussian-copula transformation) and  $F(\cdot)$  is the empirical distribution function of the bivariate data on  $\begin{bmatrix} y \\ x \end{bmatrix}$ .

$\mathbf{R}$  is the correlation matrix.

# Copula Transformed Multiple Regression: Analyses

To make these predictions independent of the model fitting process, in all cases, divide the data into training and test sets by respectively assigning the odd (even) numbered observations to the training (test) sets.

Since the Gaussian copula is used, data on all the transformed variables have zero mean and unit standard deviation.

That is however a non-issue, since,  $R^2$  as well as  $F$ -test for the model are invariant of such a transformation.

Again, since predictions are obtained in both the cases, for the data on original scale, such a location shift and scaling change do not figure in the comparison.

# Copula Transformed Multiple Regression: Algorithm

Algorithmically, the following steps are adopted in the sequence.

1. Transform the training raw data on random vector  $\begin{pmatrix} Y \\ \mathbf{x} \end{pmatrix}$  to data on uniform random variables  $U_Y, U_{X_1}, \dots, U_{X_k}$  using the empirical cumulative distribution function estimated from the data. From the estimated covariance matrix, a correlation matrix for  $\begin{pmatrix} Y \\ \mathbf{x} \end{pmatrix}$  and the corresponding empirical correlation matrix  $\Sigma$  of  $\mathbf{U} = (U_Y, U_{X_1}, \dots, U_{X_k})'$  are obtained. These provide the estimates of the copula parameters.

# Copula Transformed Multiple Regression: Algorithm

2. This is our target correlation matrix and we want, more or less, the same correlation among the multivariate uniform variables and among the multivariate normal transformed variables. Using the inverse multivariate normal cumulative distribution function on  $\mathbf{U}$ , we obtain the transformed data which are jointly distributed as the multivariate normal. We denote this by  $\begin{pmatrix} Y^* \\ \mathbf{X}^* \end{pmatrix}$ .

# Copula Transformed Multiple Regression: Algorithm

3. We fit separately, the two models which are functionally similar, yet for the above two different data sets. Specifically, these are (\* indicates the copula-transformed data),

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k + \epsilon \quad (6)$$

and

$$y^* = \beta_0^* + \beta_1^* x_1^* + \beta_2^* x_2^* + \cdots + \beta_k^* x_k^* + \epsilon^*. \quad (7)$$

# Copula Transformed Multiple Regression: Algorithm

4. For predictions, comparison is appropriate only in the original scale. That is readily available for the first model.

However, the second model is fit for transformed data and so original scale is lost. This is, however, of little concern as the correspondence between the observations in the two data sets is one to one, and thus for the test data, predicted values of  $y$  can be obtained through its correspondence with  $y^*$ .

# Copula Transformed Multiple Regression: Algorithm

Let the corresponding two predicted values of  $y$  be  $\hat{y}$  and  $\hat{y}_c$  respectively. Then for the test data set, the two sum of squared prediction errors (SSPE) are given by

$$SSPE_{Raw} = \sum_{test\ data} (y_i - \hat{y}_i)^2 \quad (8)$$

and

$$SSPE_{Copula} = \sum_{test\ data} (y_i - \hat{y}_{c,i})^2. \quad (9)$$

# Copula Transformed Multiple Regression: Algorithm

To obtain the prediction interval (say for a future observation) on the original scale and using model in (??), a little more care is needed. For a given  $\mathbf{X} = \mathbf{x}_f$  (and hence  $\mathbf{X}^* = \mathbf{x}_f^*$ ), denote the predicted value of  $y^*$  using model (??) by  $\hat{y}_f^*$  and let the corresponding prediction interval be  $(\hat{y}_{f,L}^*, \hat{y}_{f,U}^*)$ . Since  $\hat{y}_f^*$ ,  $\hat{y}_{f,L}^*$  and  $\hat{y}_{f,U}^*$  are the quantities about a future incoming observation, correspondence may not be readily available within the data set.



# Copula Transformed Multiple Regression: Wicklin's Data

We circumvent this problem by simulating a large number of observations from the  $k$ -dimensional multivariate uniform distribution corresponding to our copula, and compute the corresponding values of  $y$  and  $y^*$ . Let these simulated quantities be denoted by placing a tilde ( $\sim$ ) above the corresponding variable.

If  $\widehat{y}_f^*$  is sandwiched between two such (closest) simulated values, say  $\widetilde{y}_t^*$  and  $\widetilde{y}_{t+1}^*$ , then the predicted value of  $y$ , say  $\widehat{y}_{c,f}$  can be obtained by interpolation from  $\widetilde{y}_t$  and  $\widetilde{y}_{t+1}$ , each of which has one to one correspondence with  $\widetilde{y}_t^*$  and  $\widetilde{y}_{t+1}^*$  via  $\widetilde{u}_{y,t}$  and  $\widetilde{u}_{y,t+1}$ .

The same procedure is followed to interpolate the two prediction limits corresponding to  $\widehat{y}_{f,L}^*$  and  $\widehat{y}_{f,U}^*$ . Accordingly, a prediction interval  $(\widehat{y}_{c,f,L}, \widehat{y}_{c,f,U})$  is obtained.

# Copula Transformed Multiple Regression: Wicklin's Data

We have done this for our data sets for all the observations and plotted them against the serial number, which represents the increasing order of the (raw) data, on the response variables.

Note that this approach will also be applicable and should be followed, in the real situations when the prediction is an important objective.

# Copula Transformed Multiple Regression: Wicklin's Data

I will describe the analysis of Wicklin's data in detail so as to fully appreciate the steps of the modeling and interpretation.

The data set, consisting of 100 observations, is first arranged in the increasing order of the response variable.

We have divided the data into training data and test data, each consisting of fifty observations. Increasing order of values on response variable and taking alternative values in the training data and test data, respectively, ensure, that the two data sets are largely similar and represent the same underlying population.

$R^2$ , Adjusted  $R^2$ ,  $p$ -values and  $F$ -tests corresponding to model are obtained for the training data. For prediction, test data will be used.

# Copula Transformed Multiple Regression: Wicklin's Data

## **Residual Plots:**

Figures 2 and 3 respectively represent the residual plots for the training data for the traditional multiple regression analysis of raw data and that for the Gaussian copula transformed data. The patterns in Figure 2 clearly indicate non-randomness and a poor fit. On the contrary, the residual plot is near-ideal in Figure 3.

# Copula Transformed Multiple Regression: Wicklin's Data

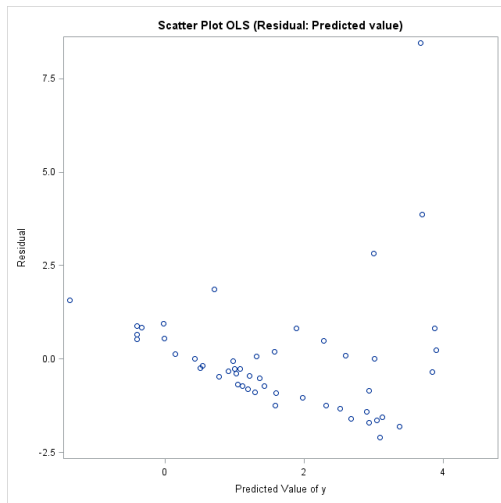


Figure: Figure 2: Wicklin's data: Scatter plot of residuals for raw training

# Copula Transformed Multiple Regression: Wicklin's Data

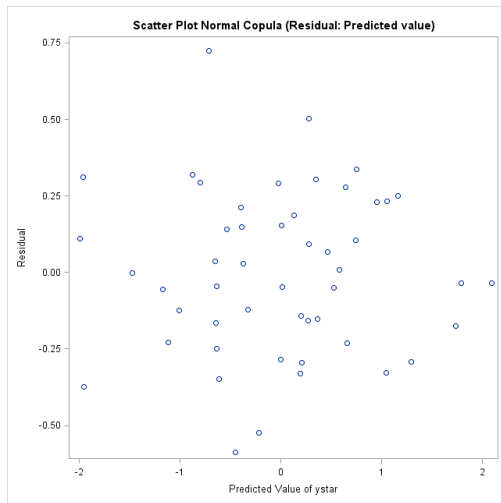


Figure: Figure 3: Wicklin's data: Scatter plot of residuals for copula

# Copula Transformed Multiple Regression: Wicklin's Data

## **QQ Plots:**

The same contrast between the two approaches is found between the two  $Q - Q$  plots, given in Figures 4 and 5 for the corresponding residuals.

# Copula Transformed Multiple Regression: Wicklin's Data

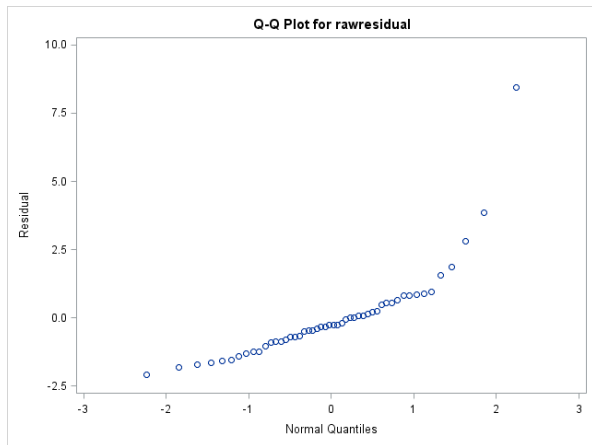


Figure 4: Wicklin's data: Residual  $Q - Q$  plot for the raw training data



# Copula Transformed Multiple Regression: Wicklin's Data

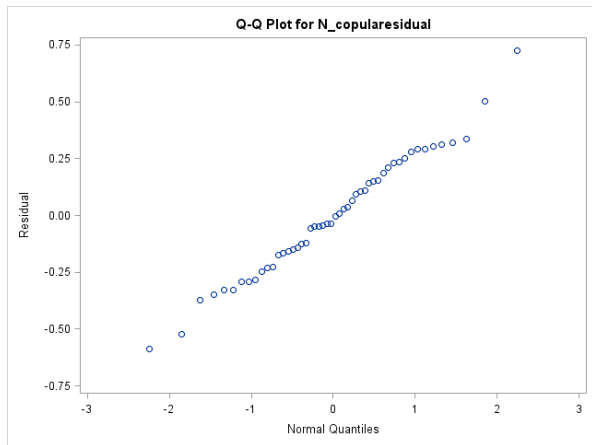


Figure 5: Wicklin's data: Residual  $Q - Q$  plot for the copula training data

# Copula Transformed Multiple Regression: Wicklin's Data

Table 1 gives the values corresponding to model fit and the statistical significance of the model. Drastic improvement in  $R^2$  (38.86% vs. 92.47%) and Adjusted  $R^2$  values is established. The same can be said about model  $F$ -statistics and corresponding  $p$ -values.

# Copula Transformed Multiple Regression: Wicklin's Data

Table: Table 1: Model fit and model significance Statistics: Wicklin's Data ( $n = 50$  for Training Data)

Statistics	Model Based on	
	Raw Data	Gaussian Copula Transformed Data
$R^2$	0.389	0.923
Adjusted $R^2$	0.349	0.920
Model		
$F$ -Stat $df$ (3,46)	9.740	188.280
$p$ -Value	$< 0.0001$	$\lll 0.0001$
Average Squared Prediction Error (Test Data $n = 50$ )	6.997	3.383

# Copula Transformed Multiple Regression: Wicklin's Data

Table 2, shows point predictions for the data along with corresponding 95% prediction intervals.

For the sake of brevity, only the first ten, last ten and middle ten observations of test data are presented. The superiority of Gaussian copula based approach is readily seen. All the point predictions using this approach are closer to the true observed responses. For the observation number 50 of the table, the observed value of the response variable is relatively very large.

# Copula Transformed Multiple Regression: Wicklin's Data

Table: Table 2A: Wicklin's data: Comparison between raw and copula regression models for test data

Obs.	$y$	$\hat{y}$	$\hat{y}^*$	95% Raw Pred. Int.	95% Copula Pred. Int.
1	0.140	-1.265	0.144	( -5.050 , 2.521 )	( 0.132 , 0.271 )
2	0.248	-0.359	0.248	( -4.076 , 3.359 )	( 0.133 , 0.356 )
3	0.249	0.085	0.284	( -3.496 , 3.666 )	( 0.168 , 0.396 )
4	0.283	0.063	0.334	( -3.565 , 3.691 )	( 0.247 , 0.492 )
5	0.293	0.148	0.284	( -3.440 , 3.736 )	( 0.168 , 0.396 )
6	0.316	0.539	0.374	( -2.998 , 4.076 )	( 0.254 , 0.615 )
7	0.353	1.015	0.360	( -2.553 , 4.583 )	( 0.249 , 0.570 )
8	0.369	-0.913	0.306	( -4.703 , 2.876 )	( 0.222 , 0.484 )
9	0.385	1.536	0.485	( -2.015 , 5.088 )	( 0.321 , 0.913 )
10	0.396	1.268	0.415	( -2.268 , 4.803 )	( 0.291 , 0.761 )

# Copula Transformed Multiple Regression: Wicklin's Data

Table: Table 2B: Wicklin's data: Comparison between raw and copula regression models for test data

Obs.	$y$	$\hat{y}$	$\hat{y}^*$	95% Raw Pred. Int.	95% Copula Pred. Int.
21	0.751	1.781	0.734	( -1.762 , 5.324 )	( 0.396 , 1.217 )
22	0.770	-0.284	0.853	( -4.034 , 3.467 )	( 0.437 , 1.567 )
23	0.822	0.455	0.923	( -3.194 , 4.105 )	( 0.487 , 1.570 )
24	0.854	1.967	0.749	( -1.547 , 5.480 )	( 0.398 , 1.350 )
25	0.922	1.731	1.079	( -1.807 , 5.269 )	( 0.669 , 2.306 )
26	0.933	2.989	0.965	( -0.583 , 6.560 )	( 0.540 , 1.771 )
27	0.981	2.526	1.414	( -1.030 , 6.082 )	( 0.801 , 2.758 )
28	1.002	2.852	0.849	( -0.745 , 6.448 )	( 0.451 , 1.498 )
29	1.078	1.230	0.918	( -2.321 , 4.780 )	( 0.487 , 1.564 )
30	1.133	2.267	1.207	( -1.338 , 5.872 )	( 0.701 , 2.686 )

# Copula Transformed Multiple Regression: Wicklin's Data

Table: Table 2C: Wicklin's data: Comparison between raw and copula regression models for test data

Obs.	$y$	$\hat{y}$	$\hat{y}^*$	95% Raw Pred. Int.	95% Copula Pred. Int.
41	2.689	1.925	2.685	( -1.627 , 5.478 )	( 1.212 , 4.758 )
42	2.751	3.061	3.191	( -0.489 , 6.612 )	( 1.563 , 7.035 )
43	2.763	3.589	2.752	( -0.007 , 7.186 )	( 1.397 , 5.562 )
44	3.456	1.830	4.112	( -1.726 , 5.387 )	( 1.865 , 8.987 )
45	3.951	4.698	4.625	( 0.914 , 8.481 )	( 2.559 , 11.435 )
46	4.290	3.753	2.758	( 0.103 , 7.402 )	( 1.400 , 5.999 )
47	4.883	3.236	4.860	( -0.383 , 6.855 )	( 2.690 , 12.911 )
48	6.952	3.790	5.731	( 0.068 , 7.513 )	( 2.747 , 16.067 )
49	9.133	3.052	9.577	( -0.499 , 6.604 )	( 4.236 , 20.029 )
50	20.029	4.557	7.294	( 0.809 , 8.305 )	( 3.426 , 18.878 )

# Copula Transformed Multiple Regression: Wicklin's Data

For Obs. 50: Observed response is very large; Both approaches underpredict the true response. Yet, the Gaussian-copula based prediction is still closer.

The prediction intervals as given in Table 2 and also graphed in Figures 6 and 7, further show that the prediction intervals are usually (and considerably) narrower when the approach is based on Gaussian copula, as compared to the raw data based regression.

The only few exceptions occur for the later few observations, but as seen in Table 2 for the last two observations, the prediction intervals based on usual regression analysis of data do not even contain the true observed value, while those based on the copula-regression approach do so for all observations except the last one.



# Copula Transformed Multiple Regression: Wicklin's Data

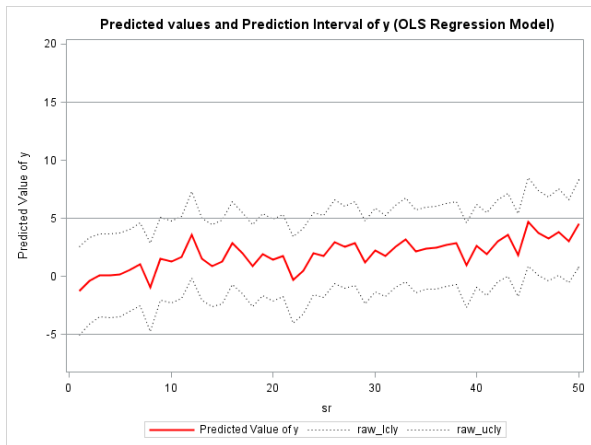


Figure 6: Wicklin's data: predicted values and prediction interval graphs for the raw test data

# Copula Transformed Multiple Regression: Wicklin's Data

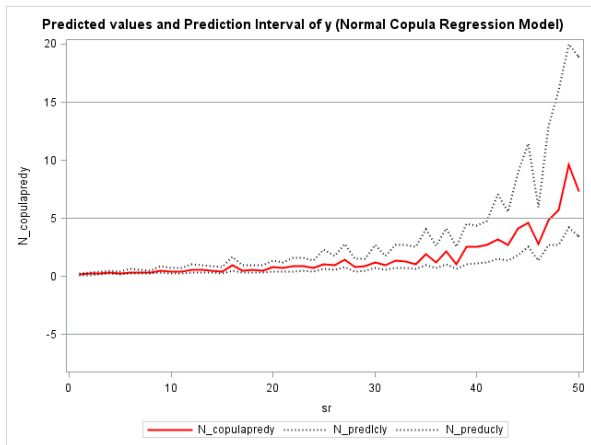


Figure 7: Wicklin's data: Predicted values and prediction interval graphs for the copula test data

# Copula Transformed Multiple Regression: Other Data Sets

The residuals from the regression for the raw data for many of the datasets indicated earlier exhibit the violation of multivariate normality and linearity of regression.

The copula transformation to multivariate normality all together circumvents these issues rather than diagnosing and correcting each of them one by one. As Cherubini, Gobbi, Mulinacci and Romagnoli in their book (p. 30), explicitly point out,

# Copula Transformed Multiple Regression: Wicklin's Data

*“... We may use the Gaussian copula whenever we want to preserve a Gaussian kind of dependence, even though the marginal distributions are not Gaussian. This dependence structure is radially symmetric and does not display tail dependence. So, we may use it if we want to link together variables whose distributions may well be assymmetric and ‘fat-tailed’, but with a dependence structure that does not change with ... .”*

In fact, this is the property which makes our approach a “general purpose” one. This simply makes the regression model realistic and inference meaningful for the transformed data and due to one-to-one correspondence with original raw data.

# Copula Transformed Multiple Regression: Wicklin's Data

**Table:** Table 3: A comparison of the two regression models for various data sets (R = Raw data; G-C = Gaussian Copula Transformed)

Sr. No.	Data	Model based on <sup>a</sup>	Est. Skewness		$R^2$	Adj. $R^2$	$p$ -Val Model	Ave Sqrd Pred Err (Test Data)
			Mardia's Skewness ( $\hat{\beta}$ )	PC Skewness ( $\hat{\eta}$ ) <sup>b</sup>				
1.	Wicklin's $n_{training} = 50$ $n_{test} = 50$	R	35.724	0.412	0.389	0.349	< 0.0001	6.997
		G-C	0.439	0.024	0.923	0.920	< 0.0001	3.383
2.	Financial Indexes $n_{training} = 508$ $n_{test} = 507$	R	1.984	0.019	0.747	0.744	< 0.0001	0.469
		G-C	0.914	0.011	0.741	0.737	< 0.0001	0.447
3.	Senic $n_{training} = 56$ $n_{test} = 56$	R	30.223	0.293	0.615	0.550	< 0.0001	1.387
		G-C	12.303	0.067	0.736	0.692	< 0.0001	1.168
4.	Prostate Cancer $n_{training} = 49$ $n_{test} = 48$	R	93.659	0.865	0.581	0.532	< 0.0001	1249.980
		G-C	8.108	0.056	0.637	0.595	< 0.0001	1290.960
5.	Real Estate Sales $n_{training} = 261$ $n_{test} = 261$	R	12.023	0.398	0.734	0.728	< 0.0001	4717109105
		G-C	2.848	0.047	0.789	0.784	< 0.0001	444560521
6.	Used Car							

# Imputing the missing values through copula transformation

Denote the fully observed variables (complete covariates) by  $\mathbf{X} = (\mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{mis}})'$  and variable with missing values by  $Y = (Y_{\text{obs}}, Y_{\text{mis}})'$ , where  $\mathbf{X}_{\text{obs}}$  and  $\mathbf{X}_{\text{mis}}$  represent the subset of  $\mathbf{X}$  for observed data  $Y_{\text{obs}}$  and missing data  $Y_{\text{mis}}$ , respectively.

Assume MCAR (Missing completely at random) scheme.

# Imputing the missing values through copula transformation

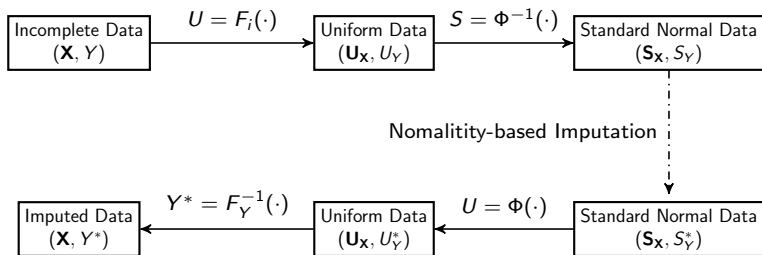


Figure: Figure 8: Procedure of imputation implementation using copula transformed data.

# Imputing the missing values through copula transformation

## Algorithm - Univariate Missing Data Pattern:

- Transform complete covariates  $\mathbf{X}$  to uniformly distributed data  $\mathbf{U}_X$  by empirical distribution.
- Transform variable  $Y$  to uniformly distributed data  $U_Y$  by empirical distribution which uses only observed data  $Y_{obs}$ .
- Convert the data  $(\mathbf{U}_X, U_Y)$  to standard normal data  $(\mathbf{S}_X, S_Y)$  using standard inverse multivariate normal cumulative distribution. That is, for each column vector  $S_i = \Phi^{-1}(U_i)$ . At this stage, the data set  $(\mathbf{S}_X, S_Y)$  are viewed as following multivariate normal distribution.



# Imputing the missing values through copula transformation

- Use one of the imputation procedures (e.g. regression, MCMC, FCS) as desired, to impute all missing values and obtain dataset  $(\mathbf{S}_X, S_Y^*)$  with imputed data.
- Back-transform the filled-in data to original scale via  $U_Y^* = \Phi(S_Y^*)$  according to the inverse of empirical marginal distribution of  $Y$ , i.e.,  $Y^* = F_Y^{-1}(U_Y^*)$ .

# Imputing the missing values through copula transformation

We apply the Iman-Conover method to generate skewed multivariate datasets.

The reason we chose this method is that we can specify the marginal distribution of each variable and also the correlation structure.

We design two groups for multivariate data setting with marginals of components as follows.

**Table:** Table 7: Marginal distributions of simulated data sets using Iman-Conover method

Group	$X_1$	$X_2$	$X_3$	$X_4$
1	Log-normal $(0, \sigma)$	Pareto $(1,1)$	Normal $(0, 1)$	Uniform $(0, 1)$
2	Log-normal $(0, \sigma)$	Normal $(0,1)$	Exp $(1)$	Uniform $(0, 1)$

# Imputing the missing values through copula transformation

We select  $X_1$  as missing variate under MCAR. In above,  $\sigma$  is set as 1.0, 2.0 and 3.0.

In each case, the following correlation structures are used.

Obviously, the data generated from above groups are clearly non-normal and  $X_1$  has larger skewness as  $\sigma$  increases.

$$\text{Corr}_1 = \begin{pmatrix} 1 & 0.75 & -0.7 & 0 \\ 0.75 & 1 & -0.95 & 0 \\ -0.7 & -0.95 & 1 & -0.2 \\ 0 & 0 & -0.2 & 1 \end{pmatrix} \quad \text{Corr}_2 = \begin{pmatrix} 1 & 0.78 & -0.67 & 0.78 \\ 0.78 & 1 & -0.89 & 0.61 \\ -0.67 & -0.89 & 1 & -0.24 \\ 0.78 & 0.61 & -0.24 & 1 \end{pmatrix}$$

$$\text{Corr}_3 = \begin{pmatrix} 1 & \rho & \rho & \rho \\ \rho & 1 & \rho & \rho \\ \rho & \rho & 1 & \rho \\ \rho & \rho & \rho & 1 \end{pmatrix} \quad \text{where } \rho \text{ is set as } 0.5, 0.6, 0.7, 0.8, \text{ and } 0.9.$$

# Imputing the missing values through copula transformation

The sample size is taken as 100 and the number of missing cases as 5.

To evaluate the quality of imputation, simulate each scenario NSIM=1,000 times and  $k$  imputation(s) and compute the mean of the sum of squared residuals by

$$\text{MSSR} = \frac{1}{\text{NSIM}} \sum_{m=1}^{\text{NSIM}} \sum_{i=1}^k \sum_{i=1}^5 \left( X_{1i}^{\text{impt}(m)} - X_{1i}^{\text{true}} \right)^2.$$

where  $X_{1i}^{\text{impt}(m)}$  is the  $m$ -th imputed value for the  $i$ -th missing value  $X_{1i}$  and  $X_{1i}^{\text{true}}$  is the true observed value of  $X_{1i}$ . Here  $k = 1$  for single imputation and  $k > 1$  for multiple imputations.

# Imputing the missing values through copula transformation

Table: Table 8: Comparison between original data and copula-transformed data using single imputation for Group 1

$\sigma$	Correlation Structure	Orig.(nor.)	MSSR Cop-tran.	Ratio ( $O/C$ )	% SSR ( $O > C$ )
1.0	Corr <sub>1</sub>	26.25	13.12	2.00	61.6
	Corr <sub>2</sub>	505.59	5.52	91.60	88.3
	Corr <sub>3</sub> ( $\rho = 0.5$ )	60.42	16.06	3.76	56.4
	Corr <sub>3</sub> ( $\rho = 0.9$ )	268.67	6.69	40.16	71.7
2.0	Corr <sub>1</sub>	6,444.07	3,625.55	1.78	71.1
	Corr <sub>2</sub>	84,351.99	2,799.91	30.13	89.0
	Corr <sub>3</sub> ( $\rho = 0.5$ )	12,519.63	3,851.05	3.25	67.2
	Corr <sub>3</sub> ( $\rho = 0.9$ )	46,791.20	2,936.43	15.93	80.1
3.0	Corr <sub>1</sub>	3,594,039.00	1,768,944.55	2.03	79.0
	Corr <sub>2</sub>	13,369,528.18	1,710,317.53	7.82	90.4
	Corr <sub>3</sub> ( $\rho = 0.5$ )	5,349,742.97	1,792,109.32	2.99	77.8
	Corr <sub>3</sub> ( $\rho = 0.9$ )	10,570,811.67	1,685,064.24	6.27	85.3

# Imputing the missing values through copula transformation

Table: Table 9: Comparison between original data and copula-transformed data using single imputation for Group 2

$\sigma$	Correlation Structure	Orig.(nor.)	MSSR Cop-tran.	Ratio ( $O/C$ )	% SSR ( $O > C$ )
1.0	Corr <sub>1</sub>	13.64	13.27	1.03	58.3
	Corr <sub>2</sub>	12.38	6.20	2.00	92.0
	Corr <sub>3</sub> ( $\rho = 0.5$ )	16.29	16.24	1.00	55.4
	Corr <sub>3</sub> ( $\rho = 0.9$ )	8.70	6.48	1.34	72.2
2.0	Corr <sub>1</sub>	4,141.31	3,614.82	1.15	73.6
	Corr <sub>2</sub>	4,220.18	3,273.75	1.29	91.2
	Corr <sub>3</sub> ( $\rho = 0.5$ )	4,264.53	3,851.76	1.11	70.2
	Corr <sub>3</sub> ( $\rho = 0.9$ )	4,197.88	2,793.05	1.50	84.1
3.0	Corr <sub>1</sub>	2,773,953.52	1,760,982.56	1.58	81.7
	Corr <sub>2</sub>	2,873,573.64	1,885,597.04	1.52	91.0
	Corr <sub>3</sub> ( $\rho = 0.5$ )	2,603,019.87	1,790,836.43	1.45	80.1
	Corr <sub>3</sub> ( $\rho = 0.9$ )	3,436,021.98	1,604,055.17	2.14	88.4

# Imputing the missing values through copula transformation

**Table:** Table 10: Comparison between original data and copula-transformed data using FCS Regression multiple imputation for **Group 1**

$\sigma$	Correlation Structure	MSSR			% SSR ( $O > C$ )
		Orig.(nor)	Cop.-tran.	Ratio ( $O/C$ )	
1.0	$\text{Corr}_3(\rho = 0.5)$	405.76	125.18	3.24	85.40%
	$\text{Corr}_3(\rho = 0.6)$	310.23	112.44	2.76	84.10%
	$\text{Corr}_3(\rho = 0.7)$	314.34	97.80	3.21	83.90%
	$\text{Corr}_3(\rho = 0.8)$	3,195.58	77.38	41.30	85.20%
	$\text{Corr}_3(\rho = 0.9)$	1,373.60	51.55	26.65	87.90%
2.0	$\text{Corr}_3(\rho = 0.5)$	104,406.86	42,081.38	2.48	85.20%
	$\text{Corr}_3(\rho = 0.6)$	193,498.66	40,986.46	4.72	85.50%
	$\text{Corr}_3(\rho = 0.7)$	210,846.48	40,019.30	5.27	85.30%
	$\text{Corr}_3(\rho = 0.8)$	491,084.91	36,171.96	13.58	86.70%
	$\text{Corr}_3(\rho = 0.9)$	260,022.58	28,639.43	9.08	89.20%
3.0	$\text{Corr}_3(\rho = 0.5)$	96,342,688.72	46,027,573.53	2.09	86.20%
	$\text{Corr}_3(\rho = 0.6)$	224,245,845.70	45,842,238.39	4.89	86.30%
	$\text{Corr}_3(\rho = 0.7)$	247,254,700.98	46,013,063.35	5.37	86.50%
	$\text{Corr}_3(\rho = 0.8)$	126,392,828.10	44,373,966.34	2.85	87.80%
	$\text{Corr}_3(\rho = 0.9)$	113,277,135.61	37,978,413.33	2.98	90.00%

# Imputing the missing values through copula transformation

**Table:** Table 11: Comparison between original data and copula-transformed data using FCS Regression multiple imputation for **Group 2**

$\sigma$	Correlation Structure	MSSR			% SSR ( $O > C$ )
		Orig.(nor)	Cop.-tran.	Ratio ( $O/C$ )	
1.0	$\text{Corr}_3(\rho = 0.5)$	164.68	123.08	1.34	85.00%
	$\text{Corr}_3(\rho = 0.6)$	149.42	109.95	1.36	83.00%
	$\text{Corr}_3(\rho = 0.7)$	132.64	95.48	1.39	82.20%
	$\text{Corr}_3(\rho = 0.8)$	112.58	79.57	1.41	83.10%
	$\text{Corr}_3(\rho = 0.9)$	84.70	53.40	1.59	86.40%
2.0	$\text{Corr}_3(\rho = 0.5)$	56,391.21	39,703.39	1.42	86.80%
	$\text{Corr}_3(\rho = 0.6)$	56,417.01	37,341.96	1.51	86.30%
	$\text{Corr}_3(\rho = 0.7)$	55,458.96	35,888.29	1.55	86.60%
	$\text{Corr}_3(\rho = 0.8)$	53,648.33	39,250.03	1.37	86.50%
	$\text{Corr}_3(\rho = 0.9)$	50,264.93	30,896.51	1.63	90.30%
3.0	$\text{Corr}_3(\rho = 0.5)$	74,206,031.68	44,532,451.31	1.67	86.90%
	$\text{Corr}_3(\rho = 0.6)$	77,116,532.29	42,631,176.04	1.81	87.50%
	$\text{Corr}_3(\rho = 0.7)$	76,957,607.66	41,943,270.76	1.83	88.00%
	$\text{Corr}_3(\rho = 0.8)$	75,842,964.49	55,855,584.64	1.36	87.70%
	$\text{Corr}_3(\rho = 0.9)$	75,184,163.76	42,523,936.26	1.77	90.80%



# Imputing the missing values through copula transformation

**Table:** Table 12: Comparison between original data and copula-transformed data using MCMC multiple imputation for **Group 1**

$\sigma$	Correlation Structure	MSSR			% SSR ( $O > C$ )
		Orig.(nor)	Cop.-tran.	Ratio ( $O/C$ )	
1.0	$\text{Corr}_3(\rho = 0.5)$	300.95	132.50	2.27	60.00%
	$\text{Corr}_3(\rho = 0.6)$	264.60	113.77	2.33	64.70%
	$\text{Corr}_3(\rho = 0.7)$	276.81	95.13	2.91	68.40%
	$\text{Corr}_3(\rho = 0.8)$	211.99	72.71	2.92	74.10%
	$\text{Corr}_3(\rho = 0.9)$	450.31	47.32	9.52	80.50%
2.0	$\text{Corr}_3(\rho = 0.5)$	71,610.82	36,851.77	1.94	66.10%
	$\text{Corr}_3(\rho = 0.6)$	156,457.61	34,240.23	4.57	70.20%
	$\text{Corr}_3(\rho = 0.7)$	171,850.60	32,943.64	5.22	75.00%
	$\text{Corr}_3(\rho = 0.8)$	60,038.31	26,745.73	2.24	79.40%
	$\text{Corr}_3(\rho = 0.9)$	131,115.17	21,115.43	6.21	83.90%
3.0	$\text{Corr}_3(\rho = 0.5)$	51,837,663.26	23,896,573.91	2.17	70.10%
	$\text{Corr}_3(\rho = 0.6)$	165,243,871.10	23,148,399.65	7.14	73.80%
	$\text{Corr}_3(\rho = 0.7)$	183,802,007.11	23,891,866.47	7.69	77.10%
	$\text{Corr}_3(\rho = 0.8)$	48,189,008.89	16,717,322.02	2.88	80.60%
	$\text{Corr}_3(\rho = 0.9)$	67,349,495.86	13,494,377.61	4.99	85.50%

# Imputing the missing values through copula transformation

**Table:** Table 13: Comparison between original data and copula-transformed data using MCMC multiple imputation for **Group 2**

$\sigma$	Correlation Structure	MSSR			% SSR ( $O > C$ )
		Orig.(nor)	Cop.-tran.	Ratio ( $O/C$ )	
1.0	$\text{Corr}_3(\rho = 0.5)$	127.54	132.10	0.97	57.60%
	$\text{Corr}_3(\rho = 0.6)$	115.44	113.55	1.02	61.10%
	$\text{Corr}_3(\rho = 0.7)$	102.34	94.08	1.09	65.30%
	$\text{Corr}_3(\rho = 0.8)$	86.99	74.35	1.17	72.20%
	$\text{Corr}_3(\rho = 0.9)$	65.26	48.47	1.35	81.20%
2.0	$\text{Corr}_3(\rho = 0.5)$	39,757.02	35,230.30	1.13	67.10%
	$\text{Corr}_3(\rho = 0.6)$	39,631.01	30,936.11	1.28	69.00%
	$\text{Corr}_3(\rho = 0.7)$	38,944.15	27,789.97	1.40	73.10%
	$\text{Corr}_3(\rho = 0.8)$	38,106.04	29,270.04	1.30	79.20%
	$\text{Corr}_3(\rho = 0.9)$	35,414.91	25,431.77	1.39	84.30%
3.0	$\text{Corr}_3(\rho = 0.5)$	43,374,578.45	22,551,396.98	1.92	71.60%
	$\text{Corr}_3(\rho = 0.6)$	44,902,158.06	19,198,351.45	2.34	72.20%
	$\text{Corr}_3(\rho = 0.7)$	44,977,065.21	17,619,423.92	2.55	77.10%
	$\text{Corr}_3(\rho = 0.8)$	45,876,052.72	30,256,223.22	1.52	80.60%
	$\text{Corr}_3(\rho = 0.9)$	44,235,290.84	29,549,312.57	1.50	85.50%

# Conclusions

- Much of the dependence based multivariate analyses for nonnormal data can be done using the copula transformation. However information about marginals is lost.
- Similar work has been done for principal component analyses, factor analyses, structural equation modeling.
- For missing data imputation for nonnormal situations, this approach is very handy. Further extensive studies showed that (for multivariate Lomax distribution) results under copula transformation are as good as those obtained by imputation by conditional expectations (assuming MCAR).
- Comparison of imputation done by our transformation to normality and that done by Box-Cox transformation showed our approach is much superior.