

POWER TRANSFORMATIONS WHEN FITTING  
THEORETICAL MODELS TO DATA

Raymond J. Carroll<sup>1</sup>

and

David Ruppert<sup>2</sup>

Running Title: Fitting Theoretical Models

<sup>1</sup> National Heart, Lung and Blood Institute and the University of North Carolina. Supported by the Air Force Office of Scientific Research Grant AFOSR-80-0080.

<sup>2</sup> University of North Carolina. Supported by National Science Foundation Grant MCS 8100748.

Some key words: Transformations, Box-Cox models, theoretical models, robustness.

AMS 1970 Subject Classifications: Primary 62F20; Secondary 62G35.

A B S T R A C T

We investigate power transformations in non-linear regression problems when there is a physical model for the response but little understanding of the underlying error structure. In such circumstances and unlike the ordinary power transformation model, both the response and the model must be transformed simultaneously and in the same way. We show by an asymptotic theory and a small Monte-Carlo study that for estimating the model parameters there is little cost for not knowing the correct transform a priori; this is in dramatic contrast to the results for the usual case that only the response is transformed.

## 1: INTRODUCTION

Often in scientific work, one observes data  $y$  and  $x^t = (x_1 \dots x_p)$  and postulates that these data follow a model

$$(1.1) \quad y_i = f(x_i, \theta_0), \quad i = 1, \dots, N,$$

where  $\theta_0$  is a  $k$ -parameter vector. The function  $f$  may be derived, for example, from differential equations believed to govern the physical system which gave rise to the data. The deterministic model (1.1) is often inadequate since the data exhibit random variation, but whereas  $f$  was derived from theoretical considerations, there is really no firm understanding of the mechanism producing the randomness. In this case, one typically assumes that

$$(1.2) \quad y_i = f(x_i, \theta_0) + \epsilon_i,$$

where the  $\{\epsilon_i\}$  are i.i.d.  $N(0, \sigma_0^2)$ . In those cases in which the data suggest that model (1.2) is also unsatisfactory, one might then assume that the errors are multiplicative and log-normal, so that

$$(1.3) \quad \log(y_i) = \log(f(x_i, \theta_0)) + \epsilon_i.$$

The point here is that model (1.1) is equivalent to the model

$$h(y_i) = h(f(x_i, \theta_0))$$

whenever  $h(\cdot)$  is a monotonic transformation. Therefore (1.2) and (1.3) are based on the same theoretical model, but they allow variability into the model in different fashions.

A more flexible approach is to take a sufficiently rich family of strictly monotonic transformations  $h(y, \lambda)$ , indexed by the  $m$ -vector parameter  $\lambda$ , and to assume that for some value  $\lambda_0$ .

$$(1.4) \quad h(y_i, \lambda_0) = h(f(x_i, \theta_0), \lambda_0) + \epsilon_i.$$

The model (1.4) is in the spirit of Box and Cox (1964), who suggested the family of power transformations with  $m = 1$  and

$$(1.4b) \quad \begin{aligned} h(y, \lambda) = y^{(\lambda)} &= (y^\lambda - 1)/\lambda && \text{if } \lambda \neq 0 \\ &= \log(y) && \text{if } \lambda = 0. \end{aligned}$$

However, as we will make clear, our proposed model (1.4) has greatly different ramifications than usually associated with the power family. Box and Cox (1964) used their family in a study of the transformation model

$$(1.5) \quad h(y, \lambda_0) = x^t \theta_0 + \epsilon.$$

Notice here that, unlike (1.4), the regression function in (1.5) is *not* transformed. Box and Cox sought a transformation which achieves 1) a simple, additive or linear model, 2) homoscedastic errors and 3) normally distributed errors. Our model is different. Theoretical considerations already provide a regression function. We hope to transform the response *and* the regression function simultaneously to obtain homoscedasticity and normality.

There are two reasons for using model (1.4) instead of simply fitting (1.1) by least squares or some other method. First, estimation of  $\theta_0$  based on model (1.4) should be more efficient than other methods. Second, it may

be necessary to estimate the entire conditional distribution of  $y$  given  $x$ ; if the data clearly suggest that the distributions of  $\{y_i, f(x_i, \theta_0)\}$  are not constant, one must go beyond standard regression methodology.

An example, which partly motivated the research of this paper, concerns the relationship between egg production in a fish stock and subsequent recruitment into the stock. At least for some species, as egg production increases, the change in the skewness and variance of recruitment is as large as the change in the median recruitment, and this change in distributional shape may have important implications for management of the fishery.

The outline of the paper is as follows. Section 2 discusses a current controversy concerning the model of Box and Cox. Bickel and Doksum (1981) have shown that, in model (1.5), the ML estimate of  $\theta_0$  can be much more variable when  $\lambda_0$  is estimated compared to when  $\lambda_0$  is known. In Section 3, we demonstrate for our model (1.5) an entirely different result: the ML estimate of  $\theta_0$  in model (1.4) turns out to be only slightly more variable when  $\lambda_0$  is unknown compared to when  $\lambda_0$  is known. In Section 4 we prove a considerably stronger result. By examining a weighted least absolute deviations estimator, we provide a lower bound of  $2/\pi$  on the asymptotic relative efficiency of the ML estimator of  $\theta_0$  in model (1.4) when  $\lambda_0$  is unknown compared to the MLE when  $\lambda_0$  is known.

## 2: RECENT STUDIES OF THE BOX AND COX MODEL

In Section 7 of Box and Cox's original paper they discuss the analysis of effects after transformation. They state that, after finding  $\hat{\lambda}$ , one should estimate effects (regression parameters) on the scale  $\hat{\lambda}$  which has been chosen for analysis and not on the true but unknown  $\lambda_0$  scale. However, in discussing interactions, they go on to state that "The general conclusion will be that to allow for the effect of analysing in terms of  $\hat{\lambda}$  rather than  $\lambda_0$ , the

residual degrees of freedom need only be reduced by ... the number of component parameters in  $\lambda$ ".

Box and Tiao (1968) agree, stating that the only practical effect between using  $\hat{\lambda}$  in the posterior distribution of  $\theta_0$ , rather than the true  $\lambda_0$ , is an adjustment in the degrees of freedom.

Bickel and Doksum (1981) disagree with this conclusion. Following calculations for the location problem done by Hinkley (1975) and suggestive Monte-Carlo results of Spitzer (1978) and Carroll (1980), they calculated for general regression the large sample information matrix of  $\lambda_0$ ,  $\sigma_0^2$  and  $\theta_0$ . They found that the large sample variance of  $\hat{\theta}$  is larger, often much larger, when  $\lambda_0$  is estimated compared to when  $\lambda_0$  is known. They also state that the conclusion of Box and Tiao is not correct. On a technical level, part of the discrepancy between Bickel and Doksum's and Box and Tiao's results may be due to the use of different transformations. Bickel and Doksum use (1.4b), while Box and Tiao use

$$z^{(\lambda)} = y^{(\lambda)} / (\dot{y})^{\lambda-1},$$

where  $\dot{y}$  is the geometric mean of the  $\{y_i\}$ . However, Hinkley and Runger (1982) found  $z^{(\lambda)}$  unsatisfactory in several respects. The differences may also be contextual; at the null hypothesis of no interaction effects, one *can* act as if  $\lambda_0$  were known, with an appropriate change in the degrees of freedom. See Carroll (1982) and Doksum and Wong (1981).

Since power transformations have been used often and with real satisfaction by applied statisticians, the findings of Bickel and Doksum were surprising and led to further research. Hinkley and Runger argue that the parameter  $\theta_0$  in (1.5) is not physically meaningful; it is defined in an unknown scale  $\lambda_0$  so that a unit change in  $x$  is not easily interpreted by  $\theta_0$  alone. Instead, they argue that in practice, the relevant distribution is the conditional distribution of  $\hat{\theta}$  given  $\hat{\lambda}$ . As  $N \rightarrow \infty$ , the conditional variance of  $\hat{\theta}$  given  $\hat{\lambda}$

and the variance of  $\hat{\theta}$  when  $\lambda_0$  is known converge to the same matrix. They then argue that, when analyzing  $\theta_0$ , no adjustment need be made for the fact that  $\lambda_0$  was estimated. This appealing behavior is somewhat counter-balanced by difficulties with the conditional mean in hypothesis testing in unbalanced designs, as pointed out by Carroll (1982).

Carroll and Ruppert (1981) also noticed the difficulty with interpreting  $\theta_0$  and studied predicting the median of  $y$  on the *original* data scale by backtransforming  $x^{t\hat{\theta}}$ . This idea of looking at the response surface avoids the problems of definition inherent with  $\theta_0$  being defined in an unknown or data dependent scale. They found that when predicting the median of  $y$ , the effect of not knowing  $\lambda_0$  can be large but is in general similar to the effect of adding one more regression parameter, and it is certainly much less severe than the effect when estimating  $\theta_0$ .

The above discussion establishes the extent of the controversy surrounding the Box and Cox model applied (1.5). We believe (1.4) entirely avoids this controversy. First, the parameter  $\theta_0$  has physical meaning even if  $\lambda_0$  is unknown, since  $f(x_i, \theta_0)$  is the median of  $y_i$  no matter what the true scale. Secondly, the large sample analysis to follow indicates that  $\hat{\theta}$  is only slightly more variable when  $\lambda_0$  is estimated than when  $\lambda_0$  is known.

### 3: LIKELIHOOD ANALYSIS

The likelihood analysis proceeds as follows: define

$$z_i = dh(f_i(\theta_0), \lambda_0) / d\theta_0$$

$$f_i(\theta) = f(x_i, \theta), \quad f_i = f_i(\theta_0),$$

$$h_y(y) = h_y(y, \lambda) = dh(y, \lambda) / dy, \quad \text{and } h(y) = h(y, \lambda).$$

Let  $h_{\lambda}(y)$  and  $h_{\lambda\lambda}(y)$  be the gradient vector and Hessian of  $h(y, \lambda)$  with respect to  $\lambda$ . By simple algebra we find the joint information matrix of  $(\theta_0, \sigma_0, \lambda_0)$

as (all summations are from 1 to N)

$$(3.1) \quad N^{-1} I = \begin{pmatrix} S/\sigma_o^2 & 0 & C_1/\sigma_o^2 \\ \cdot & 1/(2\sigma_o^4) & C_2/\sigma_o^4 \\ \cdot & \cdot & C_3/\sigma_o^2 \end{pmatrix}$$

where

$$(3.2) \quad \begin{aligned} S &= N^{-1} \sum z_i z_i^t \\ C_1 &= -N^{-1} E \sum z_i [h_\lambda(y_i) - h_\lambda(f_i)]^t \\ C_2 &= -N^{-1} E \sum \epsilon_i [h_\lambda(y_i) - h_\lambda(f_i)]^t \\ C_3 &= N^{-1} E \sum \{ [h_\lambda(y_i) - h_\lambda(f_i)] [h_\lambda(y_i) - h_\lambda(f_i)]^t \\ &\quad + \epsilon_i [h_{\lambda\lambda}(y_i) - h_{\lambda\lambda}(f_i)] + (\partial/\partial\lambda)(\partial/\partial\lambda)^t \log [h_y(y_i)] \}. \end{aligned}$$

Using the work of Hoadley (1971), it is straightforward, though perhaps somewhat tedious, to establish conditions sufficient that  $(\hat{\theta}, \hat{\sigma}^2, \hat{\lambda})$  is consistent and asymptotically normal. We will not pursue this matter further, but rather we will assume that  $(\hat{\theta}, \hat{\sigma}^2, \hat{\lambda})^t$  is approximately  $N((\theta_o, \sigma_o^2, \lambda_o)^t, I^{-1})$  and we will study  $I^{-1}$ .

In general,  $C_1$  and  $C_2$  are not zero and the asymptotic distribution of  $(\hat{\lambda}, \hat{\sigma}^2)$  when  $\lambda_o$  is estimated differs from when  $\lambda_o$  is known. At least to this point then, the analysis is similar to those done in the usual Box-Cox model (1.5). The key question, of course, is whether or not  $C_1$  and  $C_2$  are sufficiently different from zero to seriously affect the distribution of  $\hat{\lambda}$ .

The expressions  $C_1$ ,  $C_2$  and  $C_3$  are complex even when  $f_i(\theta_o)$  has a nice form such as simple linear regression. To simplify matters sufficiently that we can gain some insight about the difference between knowing and estima-



ting  $\lambda_0$ , we follow Bickel and Doksum and others and let  $\sigma_0 \rightarrow 0$ . While Bickel and Doksum let  $N \rightarrow \infty$  and  $\sigma_0 \rightarrow 0$  simultaneously, we let  $N \rightarrow \infty$  and then  $\sigma_0 \rightarrow 0$ . There is no essential difference between the two approaches. Our is very suitable for heuristic arguments.

It should be emphasized that we are not concerned only, or even primarily, with small  $\sigma_0$ . In fact, the need for transformation is greater when  $\sigma_0$  is large. The small  $\sigma_0$  asymptotics do, however, lead to major simplifications, and the Monte-Carlo results presented later agree with them.

Taylor expansions show that under mild regularity conditions

$$(3.3) \quad C_1 = o(\sigma_0^2), C_2 = o(\sigma_0^2), \text{ and } C_3 = o(\sigma_0^2) \text{ as } \sigma_0 \rightarrow 0.$$

Standard calculations show that when  $\lambda_0$  is known,

$$(3.4) \quad N^{1/2} \text{ Covariance } [(\hat{\theta} - \theta_0)/\sigma_0, (\hat{\sigma}^2 - \sigma_0^2)/\sigma_0^2 | \lambda_0 \text{ known}] \\ \rightarrow A^{-1} = \begin{pmatrix} S^{-1} & 0 \\ 0 & 2 \end{pmatrix}.$$

Let  $D = \text{Diag}(\sigma_0, \sigma_0^2, 1)$ . Then, to find this limiting covariance matrix when  $\lambda_0$  is unknown, we must find the upper left  $(k+1) \times (k+1)$  corner of

$$DID = \begin{pmatrix} S & 0 & C_1/\sigma_0 \\ \cdot & 2 & C_2/\sigma_0^2 \\ \cdot & \cdot & C_3/\sigma_0^2 \end{pmatrix}$$

which by standard results on inverting partitioned matrices is

$$A^{-1} + FE^{-1}F^t$$

where  $A^{-1}$  is given in (3.4),  
 $E = C_3/\sigma_0^2 - B^t A B,$   
 $F = A^{-1}B,$

and

$$B = \begin{pmatrix} C_1/\sigma_0 \\ C_2/\sigma_0^2 \end{pmatrix}.$$

Clearly,

$$F = \begin{pmatrix} S^{-1}C_1/\sigma_0 \\ 2C_2/\sigma_0 \end{pmatrix}$$

and

$$E = C_3/\sigma_0^2 - C_1^t S^{-1} C_1/\sigma_0^2 - 2C_2^t C_2/\sigma_0^4.$$

In order to obtain simple asymptotics, we will assume that for  $\sigma_0$  fixed,  $C_1/\sigma_0^2$ ,  $C_2/\sigma_0^2$ , and  $C_3/\sigma_0^2$  converge as  $N \rightarrow \infty$ , and that these, in turn, have limits  $D_1$ ,  $D_2$ , and  $D_3$  respectively as  $\sigma_0 \rightarrow 0$ . We also assume that  $S \rightarrow S_\infty$  (positive definite) as  $N \rightarrow \infty$ . If  $D_3 - 2D_2^t D_2$  is nonsingular, then

$$\lim_{\sigma_0 \rightarrow 0} \lim_{N \rightarrow \infty} F E^{-1} F^t = \begin{pmatrix} 0 & 0 \\ 0 & W \end{pmatrix}$$

where  $W = 4 D_2^t [D_3 - D_2^t D_2]^{-1} D_2$ . Therefore

$$\lim_{\sigma_0 \rightarrow 0} \lim_{N \rightarrow \infty} (A^{-1} + F E^{-1} F^t) = \begin{pmatrix} S^{-1} & 0 \\ 0 & 2 + W \end{pmatrix}$$

**THEOREM 1.** Assume that the limits  $D_1$ ,  $D_2$ ,  $D_3$ ,  $S_\infty$  mentioned above exist and that  $D_3 - 2D_2^t D_2$  is nonsingular. As  $N \rightarrow \infty$  and then  $\sigma_0 \rightarrow 0$ , the limit distribution of  $\hat{\theta}$  is the same whether  $\lambda_0$  is known or unknown. The limit distribution of  $\hat{\sigma}$  depends on whether  $\lambda_0$  is known or unknown.

As an example consider multiple linear regression and the power transformation family, i.e.,  $h(y, \lambda)$  is given by (1.4b) and

$$h(y_i, \lambda) = x_i^t \theta_0 + \epsilon$$

where  $x_1, \dots, x_n$  are known  $k \times 1$  vectors. Also, suppose that  $\lambda_0 = 0$ , i.e., the log transformation is needed. Then  $h_y(y) = y^{\lambda-1}$ ,  $h_\lambda(y) = (\log y)^{2/2}$ , and  $h_{\lambda\lambda}(y) = (\log y)^{3/3}$ . We find that

$$A = N^{-1} \sum x_i x_i^t / (x_i^t \theta_0)^2$$

$$C_1 = -(2N)^{-1} E \sum [x_i / (x_i^t \theta_0)] \{ [\log(x_i^t \theta_0) + \epsilon_i]^2 - [\log(x_i^t \theta_0)]^2 \}$$

$$= -\sigma_0^2 (2N)^{-1} \sum x_i / (x_i^t \theta_0),$$

$$C_2 = -(2N)^{-1} E \sum \epsilon_i \{ [\log(x_i^t \theta_0) + \epsilon_i]^2 - [\log(x_i^t \theta_0)]^2 \}$$

$$= -N^{-1} \sum \log(x_i^t \theta_0) \sigma_0^2,$$

and

$$C_3 = (4N)^{-1} E \sum \{ [\log(x_i^t \theta_0) + \epsilon_i]^2 - [\log(x_i^t \theta_0)]^2 \}^2$$

$$+ (3N)^{-1} E \sum \epsilon_i \{ [\log(x_i^t \theta_0) + \epsilon_i]^3 - [\log(x_i^t \theta_0)]^3 \}$$

$$= 7/4 \sigma_0^4 + 2\sigma_0^2/N \sum [\log(x_i^t \theta_0)]^2.$$

Therefore,

$$D_1 = \lim_{N \rightarrow \infty} (2N)^{-1} \sum x_i / (x_i^t \theta_0),$$

$$D_2 = \lim_{N \rightarrow \infty} N^{-1} \sum \log(x_i^t \theta_0),$$

and

$$D_3 = 2 \lim_{N \rightarrow \infty} N^{-1} \sum [\log(x_i^t \theta_0)]^2,$$

provided the above limits exist. Thus, the  $1 \times 1$  matrix  $D_3 - 2D_2^t D_2$  is twice the limit of the variance of  $\log(x_1^t \theta_0), \dots, \log(x_N^t \theta_0)$ , and will be nonsingular except in degenerate situations.

There is thus a fundamental difference between the models (1.4) and (1.5). A small simulation study is outlined in Section 6 and helps back up Theorem 1. This result can be extended to non-normal error distributions as

well as the robust methods of Carroll (1980) and Bickel and Doksum (1981). The details are not instructive.

4: A LOWER BOUND ON THE EFFICIENCY OF THE MLE.

Let  $\hat{\theta}(\hat{\lambda})$  and  $\hat{\theta}(\hat{\lambda}_o)$  denote the ML estimator with  $\lambda_o$  estimated and known respectively. Let  $ARE(\hat{\theta}_1, \hat{\theta}_2)$  be the asymptotic relative efficiency of  $\hat{\theta}_1$  to  $\hat{\theta}_2$ . For fixed  $\sigma_o$ , it is difficult to find  $ARE(\hat{\theta}(\lambda_o), \hat{\theta}(\hat{\lambda}))$  and, in fact, this may depend on  $\theta_o$ ,  $\lambda_o$ , the  $\{x_i\}$  and the coordinate of  $\theta$  being estimated. All that can be said for certain is that this ARE is at least one and converges to one as  $\sigma_o \rightarrow 0$ . In this section we will define a weighted  $L_1$  or least absolute deviation estimator  $\hat{\theta}(W)$  and show that  $ARE(\hat{\theta}(\lambda_o), \hat{\theta}(W)) \leq \pi/2$ . Under reasonable regularity conditions, this means that  $ARE(\hat{\theta}(\lambda_o), \hat{\theta}(\hat{\lambda}))$  is bounded between one and  $\pi/2$ , in vivid contrast to the Box and Cox model (1.5) in which this last ARE can approach infinity. We first look at general weighted  $L_1$  estimators. The results stated here seem to be new and are of interest in their own right.

Let  $w_1, \dots, w_N$  be positive numbers and let  $\theta(L)$  be any point which minimizes the expression

$$\sum w_i |y_i - f_i(\hat{\theta}(L))| .$$

Under (1.4),  $f_i(\theta_o)$  is the unique median of  $y_i$ , so we can expect  $\hat{\theta}(L)$  to be consistent. The unweighted  $L_1$  estimate for linear models was studied by Ruppert and Carroll (1980). Those results suggest that

$$(4.1) \quad 0 \doteq \sum w_i \text{sign}(y_i - f_i(\hat{\theta}(L)))s_i,$$

$$s_i = df_i(\theta_0)/d\theta.$$

Define  $r_i = y_i - f_i(\theta_0)$  and let  $m_i$  be the density of  $r_i$ . By a generalization of the strong law, for example Theorem 7.1 of Carroll and Ruppert (1982) which itself generalizes Lemma 4.2 of Bickel (1975),

$$(4.2) \quad 0 \doteq \sum w_i \{ \text{sign}(y_i - f_i(\hat{\theta}(L))) - \text{sign}(r_i) \} s_i$$

$$- (E \sum w_i \{ \text{sign}(y_i - f_i(\theta)) - \text{sign}(r_i) \} s_i) | \theta = \hat{\theta}(L)).$$

Now, as  $\epsilon \rightarrow 0$ , we obtain that

$$(4.3) \quad E(\text{sign}(r_i + \epsilon) - \text{sign}(r_i)) - 2\epsilon m_i(0) \rightarrow 0.$$

Combining (4.1)-(4.3) we get to order  $o(N^{-1/2})$ ,

$$(4.4) \quad (\hat{\theta}(L) - \theta_0) \doteq \frac{1}{2} (\sum w_i m_i(0) s_i s_i^t)^{-1} \sum w_i s_i \text{sign}(r_i).$$

Now, since for model (1.4)

$$\epsilon_i = h(f_i(\theta_0) + r_i, \lambda_0) - h(f_i(\theta_0), \lambda_0),$$

we then have

$$(4.5) \quad m_i(0) = (2\pi\sigma_0^2)^{-1/2} h_y(f_i(\theta_0), \lambda_0).$$

Thus, if we chose

$$(4.6) \quad w_i = h_y(f_i(\theta_0), \lambda_0),$$

we have by (4.4)-(4.6) and the Central Limit Theorem that

$$N^{1/2}(\hat{\theta}(L) - \theta_0) / \sigma_0 \xrightarrow{L} N(0, (\pi/2) S^{-1}).$$

Now  $\hat{\theta}(L)$  is not a bona fide estimator since  $w_i$  in (4.6) requires  $\lambda_0, \theta_0$  to be known. However, if in (4.6) one plugs in any  $N^{1/2}$  consistent estimators of  $\theta_0$  and  $\lambda_0$  and calls the  $L_1$  estimate based on these new weights  $\hat{\theta}(W)$ , then using Theorem 7.1 of Carroll and Ruppert (1982), one can also show that

$$N^{1/2}(\hat{\theta}(W) - \theta_0) / \sigma_0 \xrightarrow{L} N(0, (\pi/2) S^{-1}).$$

Now, because

$$N^{1/2}(\hat{\theta}(\lambda_0) - \theta_0) / \sigma_0 \xrightarrow{L} N(0, S^{-1}).$$

it then follows that

$$(4.8) \quad \begin{aligned} \text{ARE}(\hat{\theta}(\lambda_0), \hat{\theta}(W)) &= \pi/2, \\ \text{ARE}(\hat{\theta}(\lambda_0), \hat{\theta}(\hat{\lambda})) &\leq \pi/2. \end{aligned}$$

Theorem 1 and the Monte-Carlo results to follow indicate that the upper bound in (4.8) is quite conservative. The beauty of (4.8) is that it is a bound that does not depend on  $\sigma_0$ .

The weighted  $L_1$  estimator may well be useful for example if in (1.4) one suspected that the errors  $\{\epsilon_i\}$  are not normal. It is a consistent estimator of  $\theta_0$  provided that 0 is the unique median of  $\epsilon_i$ . Symmetry of  $\epsilon_i$  is not needed.

### 5: THE K-SAMPLE PROBLEM

Our model (1.4) and Theorem 1 provide some useful insight into the  $k$ -sample problem under the formulation (1.5) of Box and Cox. In their model, for each of  $k$  populations we have

$$(5.1) \quad h(y_{ij}, \lambda_0) = \mu_j + \epsilon_{ij} \quad j = 1, \dots, k; \quad i = 1, \dots, N_j.$$

The equivalent formulation from our viewpoint is

$$(5.2) \quad \begin{aligned} h(y_{ij}, \lambda_0) &= h(\xi_j, \lambda_0) + \epsilon_{ij}, \\ \mu_j &= h(\xi_j, \lambda_0). \end{aligned}$$

Here  $\xi_j$  is the median of  $y_{ij}$  on the original scale and  $\mu_j$  is the expected value of  $y_{ij}$  in the  $\lambda_0$  scale. The results of Carroll and Ruppert imply that for estimating the  $\xi$ 's, there is little cost in not knowing  $\lambda_0$ , while for estimating the  $\mu$ 's, Bickel and Doksum show that the cost of not knowing  $\lambda_0$  can be enormous. Since

$$H_0: \mu_1 = \dots = \mu_k \quad \text{iff} \quad H_0: \xi_1 = \dots = \xi_k,$$

there should be little cost in testing for equality of means when  $\lambda_0$  is unknown. These heuristics are formally proven by Carroll (1982) and Doksum and Wong (1981).

### 6: MONTE-CARLO.

To study  $\hat{\theta}$  when  $N$  is finite and  $\sigma_0$  is not necessarily small, we undertook a small simulation of the model

$$(6.1) \quad h(y_i, \lambda_0) = h(\theta_1 + \theta_2 x_i, \lambda_0) + \sigma_0 \epsilon_i,$$

where  $h(\cdot)$  is the Box and Cox power family (1.4b). In our simulations,  $N = 50$ , the design points  $\{x_i\}$  were equally spaced on  $[-1, 1]$ , the errors were normally distributed with mean zero and variance one and  $\theta_1 = 7$ ,  $\theta_2 = 2$ .

We considered three estimators:

- 1) ML estimator,  $\lambda_0$  known (KNOWN)
- 2) ML estimator,  $\lambda_0$  unknown (MLE)
- 3) The ordinary least squares estimator (LSE) without any transformation.

The median of  $y$  is  $\theta_1 + \theta_2 x$ , so that LSE forms an especially plausible estimator of the slope  $\theta_2$  (for which it is consistent). We chose three values of  $\sigma_0$ :

$$\sigma_0 = 0.05, 0.10, \text{ and } 0.50.$$

We present results in Tables 1 and 2 for  $\lambda_0 = 0$  (log-normal data) and  $\lambda_0 = 0.25$ . There were 600 replications of the experiment for each  $(\lambda_0, \epsilon_0)$  and each estimator, all generated from a common set of random numbers. The normal random deviates were generated from the IMSL routine GGNPM. Estimation of  $(\theta_1, \theta_2)$  for each  $\lambda$  was done by the IMSL routine ZXSSQ while ZXGSN was used to estimate  $\lambda_0$ .

The results for the ML estimator with  $\lambda_0$  unknown (denoted MLE) are very encouraging. The mean square errors for MLE are quite close to those for KNOWN, the ML estimator with  $\lambda_0$  known, especially for the slope  $\theta_2$ . These results agree with our small  $\sigma$  theory and indicate the minimal cost for not knowing  $\lambda_0$ . The relative efficiencies of MLE to KNOWN are always well above the lower bound of  $2/\pi$ . To appreciate how well MLE does relative to KNOWN



(line 2 of Tables 1 and 2), it is enlightening to study Table 5 of Bickel and Doksum (1981); in their model which we call (1.5), they have ratios  $\text{MLE}(\lambda_0 \text{ estimated})/\text{KNOWN}(\lambda_0 \text{ known})$  always at least 1.5 and as large as 211, while ours never exceed 1.2.

The other valuable point learned from Table 2 is that when estimating the slope  $\theta_2$ , the ML estimator MLE with  $\lambda_0$  unknown tends to dominate the LSE, especially for larger values of  $\sigma_0$ . In other words, for our model (1.4), there is real value to transformation when it is appropriate.

R E F E R E N C E S

- BICKEL, PETER J. (1975). One step Huber estimates in the linear model. *J. Amer. Statist. Assoc.* 70, 428-433.
- BICKEL, P.J., and DOKSUM, K.A. (1981). An analysis of transformations revisited. *J. Am. Statist. Assoc.* 76, 296-311.
- BOX, GEORGE E.P. and COX, DAVID R. (1964). An analysis of transformations. *J. Roy. Statist. Soc. Ser. B* 26, 211-252.
- BOX, G.E.P., and TIAO, G.C. (1973). *Bayesian Inference in Statistical Analysis*. Reading, Mass. Addison-Wesley.
- CARROLL, R.J. (1980). A robust method for testing transformations to achieve approximate normality. *J. Roy. Statist. Soc. Series B* 42, 71-78.
- CARROLL, R.J. (1982). Tests for regression parameters in power transformation models. Tentatively accepted by the *Scand. J. Statist.*
- CARROLL, R.J., and RUPPERT, D. (1981). Prediction and the power transformation family. *Biometrika* 68, 609-617.
- CARROLL, R.J. and RUPPERT, DAVID (1982). Robust estimation in heteroscedastic linear models. To appear in *Ann. Statist.*
- DOKSUM, K.A., and WONG, C.W. (1981). *Statistical tests after transformations*. Manuscript.
- HINKLEY, D.V., and RUNGER, G. (1981). Analysis of transformed data. To appear in *J. Am. Statist. Assoc.*
- HOADLEY, B.A. (1971). Asymptotic properties of maximum likelihood estimators for the independent not identically distributed case. *Ann. Math. Statist.* 42, 1977-1991.
- RUPPERT, D. and CARROLL, R.J. (1980). Trimmed least squares estimation in the linear model. *J. Am. Statist. Assoc.* 75, 828-838.
- SPITZER, J.J. (1978). A Monte-Carlo investigation of the Box-Cox transformation in small samples. *J. Am. Statist. Assoc.* 73, 488-495.

TABLE #1

Results of the Monte-Carlo study described in the text. These results are for the INTERCEPT. The median response is linear with intercept = 7 and slope = 2.

KNOWN = ML estimate with  $\lambda$  known.  
 MLE = ML estimate with  $\lambda$  unknown.  
 LSE = ordinary least squares estimate.

	$\lambda$ 0.00			$\lambda$ 0.25		
	$\sigma$ 0.05	$\sigma$ 0.10	$\sigma$ 0.50	$\sigma$ 0.05	$\sigma$ 0.10	$\sigma$ 0.50
BIAS OF KNOWN	0.03	0.06	0.56	0.01	0.03	0.23
MSE OF KNOWN	2.41	9.67	24.87	0.90	3.59	9.04
BIAS OF MLE	0.02	0.04	0.60	0.01	0.02	0.19
$\frac{\text{MSE OF MLE}}{\text{MSE OF KNOWN}}$	1.02	1.05	1.14	1.01	1.03	1.12
MSE OF MLE - MSE OF KNOWN	0.05	0.47	3.44	0.01	0.09	1.09
S.E. OF ABOVE DIFF.	0.02	0.15	0.77	0.01	0.04	0.25
BIAS OF LSE	0.11	0.40	9.48	0.04	0.13	2.60
$\frac{\text{MSE OF MLE}}{\text{MSE OF LSE}}$	0.97	0.90	0.22	1.00	0.98	0.63
MSE OF MLE - MSE OF LSE	-0.06	-1.15	-96.62	0.00	-0.06	-6.07
S.E. OF ABOVE DIFF.	0.04	0.33	4.71	0.01	0.06	0.78

In these calculations, the mean square error (MSE) and S.E. of difference terms are multiplied by  $T^2$ . Here  $T = 10$  if  $\sigma \leq 0.10$ ,  $T = 1$  if  $\sigma = 0.50$ .

TABLE #2

Results of the Monte-Carlo study described in the text. These results are for the SLOPE. The median response is linear with intercept = 7 and slope = 2.

KNOWN = ML estimate with  $\lambda$  known.  
 MLE = ML estimate with  $\lambda$  unknown.  
 LSE = ordinary least squares estimate.

$\lambda$	0.00			0.25			
	$\sigma$	0.05	0.10	0.50	0.05	0.10	0.50
BIAS OF KNOWN		0.01	0.01	0.03	0.00	0.01	0.02
MSE OF KNOWN		7.08	28.36	72.23	2.71	10.83	27.24
BIAS OF MLE		-0.01	-0.04	-0.15	0.00	-0.02	-0.16
$\frac{\text{MSE OF MLE}}{\text{MSE OF KNOWN}}$		1.06	1.06	1.01	1.06	1.06	1.03
MSE OF MLE - MSE OF KNOWN		0.41	1.57	0.95	0.15	0.60	0.72
S.E. OF DIFF.		0.10	0.40	0.67	0.04	0.77	0.27
BIAS OF LSE		0.05	0.15	2.97	0.02	0.04	0.50
$\frac{\text{MSE OF MLE}}{\text{MSE OF LSE}}$		0.98	0.96	0.59	1.01	1.01	0.91
MSE OF MLE - MSE OF LSE		-0.16	-1.29	-50.54	0.05	0.13	-2.81
S.E. OF DIFF.		0.18	0.80	5.10	0.06	0.23	0.74

In these calculations, the mean square errors (MSE) and S.E. of difference terms are multiplied by  $T^2$ . Here,  $T = 10$  if  $\sigma \leq 0.10$ ,  $T = 1$  if  $\sigma = 0.50$ .