# 行政院國家科學委員會專題研究計畫成果報告

## 用 Gamma 過程分析單調長期追蹤資料
## Analyzing Monotone Longitudinal Data via Gamma Processes

## 一、中文摘要

在這篇文章中，我們提出一般變換路徑模型來建構長期追蹤資料。我們用修正冪變換族和修正指數變族換當例子來闡明我們的模型。最後，我們討論相對應的估計方法。

關鍵詞：長期追蹤資料、一般路徑模型、變換、不等變異數、修正冪變換、修正指數變換、故障前時間分布。

## Abstract

In this paper, the general transformed path (GTP) model with heteroscedastic errors for modeling longitudinal data is proposed. The families of modified power transformations and modified exponential transformations are used as examples to illustrate our models. The corresponding estimation method is also discussed.

**Keywords**: longitudinal data, general path model, transformation, heteroscedasticity, modified power transformation, modified exponential transformation, time-to-failure distribution.

## 二、緣由與目的

In this paper, we are mainly concerned with modeling longitudinal data. In this section, we briefly introduce the general path (GP) model (Lu and Meeker, 1993) and the transformation model (e.g., see Atkinson and Cox, 1988 and Taylor, 1998) with heteroscedastic errors to proceed the discussion.

For modeling degradation data, Lu and Meeker (1993) proposed the following GP model: For each test unit $i$ in a random sample of size $n$ from the population of all units, assume that degradation measurements are available for prespecified times $t(1)$, $t(2)$, …, $t(m(i))$, where $t(j)$ is the time of the $j$th measurement or inspection on test unit $i$ for $j = 1, 2, …, m(i)$ and $m(i)$ is the total number of inspections on test unit $i$. For each test unit $i$, let $x_{i,t} \geq 0$ be its actual path for $t \geq 0$. Then the sample path $y_{i,t(1)}, y_{i,t(2)}, …, y_{i,t(m(i))}$ of test unit $i$ for $t = t(1), t(2), …, t(m(i))$ is given by

$$h(y_{i,t(j)}) = h(x_{i,t(j)}) + V_{i,t(j)}$$
$$= \sim(t(j); S_0, S_i) + V_{i,t(j)},$$

$j = 1, 2, …, m(i)$; $i = 1, 2, …, n$, where $h$ is a known strictly monotone transformation, e.g., the log transformation, $V_{i,t(j)}$'s are i.i.d. $N(0, t^2)$ measurement errors with unknown variance $t^2 > 0$, $\sim$ is a known regression function indexed by a fixed-effects parameter vector $S_0$ and a random-effects parameter vector $S_i$, $S_i$'s are independent of $V_{i,t(j)}$'s and i.i.d. $N(0, d_S)$ with unknown covariance matrix $d_S$, and $\sim(t, S_0, S_i) = h(x_{i,t})$ for $t \geq 0$. Let $h^{-1}$ be the inverse function of $h$ and set $(h^{-1})^{(k)}(v) = \partial^k h^{-1}(v)/\partial v^k$ for $k = 1, 2, …$. Then, for each $(i,j)$ pair, $x_{i,t(j)} = h^{-1}(\sim(t(j); S_0, S_i))$ is the conditional median of $y_{i,t(j)}$ given $S_i$, but in general not the conditional mean of $y_{i,t(j)}$ given $S_i$. Since all measurement errors are normally distributed, the range of the transformation is $R$, where $R = (-\infty, \infty)$. In the literature, the most commonly used transformation for

positive data with this property is the log transformation. In such a situation, for each $(i,j)$ pair, $x_{i,t(j)} = \exp(\sim(t(j);s_0,s_i))$ is the conditional median of $y_{i,t(j)}$ given $s_i$, but less than $\exp(t^2/2)\cdot x_{i,t(j)}$, the conditional mean of $y_{i,t(j)}$ given $s_i$.

For modeling independent and continuous data, the transformation model with heteroscedastic errors is as follows:
$$h(y_i; \gamma) = \sim(x_i; s) + g(\sim(x_i; s), z_i; x)\cdot v_i,$$
$i = 1, 2, \ldots, n$, where $y_i$ is the observation for subject $i$, $h$ is a known strictly monotone transformation indexed by a transformation parameter vector $\gamma$, both $x_i$ and $z_i$ are known covariates of subject $i$, $\sim$ is a known regression function indexed by a regression parameter vector $s$, $g$ is a known positive weight function indexed by a variance parameter vector $x$, and $v_i$'s are *i.i.d.* $N(0,1)$ standardized errors. Since all standardized errors are normally distributed, the range of the transformation is $R$. For each fixed $\gamma$, let $h^{-1}(\cdot; \gamma)$ be the inverse function of $h(\cdot; \gamma)$ and set $(h^{-1})^{(k)}(v; \gamma) = \partial^k h^{-1}(v; \gamma)/\partial v^k$ for $k = 1, 2, \ldots$. Then, for each $i$, $h^{-1}(\sim(x_i; s); \gamma)$ is the median of $y_i$, but in general not the mean of $y_i$.

In the literature, the classical likelihood inference in transformation models for continuous data is usually based on the key assumption that transformed data are normally distributed, e.g., the transformation model with heteroscedastic errors. However, this assumption would fail for the case where the range of the transformation is not $R$. As an example, the family of modified power transformations (Box and Cox, 1964),
$$h(u; \gamma)$$
$$= 1_{R\backslash\{0\}}(\gamma)\cdot(u^{\gamma}-1)/\gamma + 1_{\{0\}}(\gamma)\cdot\log(u) \equiv u^{(\gamma)},$$
$u > 0$, is most commonly used to transform positive data, where $1_S(\gamma) = 1$ if $\gamma \in S$ and 0 if $\gamma \in R\backslash S$ for any subset $S$ of $R$. By applying a modified power transformation to positive data, the range of the transformation is not $R$ except for the log transformation. As another example, the family of modified exponential transformations,
$$h(u; \gamma)$$
$$= 1_{R\backslash\{0\}}(\gamma)\cdot[\exp(\gamma\cdot u)-1]/\gamma + 1_{\{0\}}(\gamma)\cdot u \equiv u^{[\gamma]},$$
$u \in R$, can be used to transform real-valued data. By applying a modified exponential transformation to real-valued data, the range of the transformation is not $R$ except for the identity transformation. Furthermore, when the range of the transformation is not $R$, the standardized errors in the transformation model with heteroscedastic errors have possibly different supports and thus are not necessarily identically distributed. One way to tackle this problem is only to assume that all standardized errors in the transformation model with heteroscedastic errors have mean 0 and variance 1, but not necessarily normally and/or identically distributed.

## 三、結果與討論

In this section, we propose the GTP model with heteroscedastic errors for modeling longitudinal data as follows: For each subject $i$ in a random sample of size $n$ from the population of all possible subjects, assume that repeated measurements are available for prespecified times $t(i,1)$, $t(i,2)$, …, $t(i,m(i))$, where $t(i,j)$ is the time of the $j$th measurement of subject $i$ for $j = 1, 2, \ldots, m(i)$ and $m(i)$ is the total number of repeated measurements of subject $i$. For each subject $i$, let $\{x_{i,t}: t \geq 0\}$ be its actual path for $t \geq 0$. Then the sample path $y_{i,t(i,1)}$, $y_{i,t(i,2)}$, …, $y_{i,t(i,m(i))}$ of subject $i$ for $t = t(i,1)$, $t(i,2)$, …, $t(i,m(i))$ is given by
$$h(y_{i,t(i,j)}; \gamma_0, \gamma_i)$$
$$= h(x_{i,t(i,j)}; \gamma_0, \gamma_i) + v_{i,t(i,j)}^*$$
$$= \sim(t(i,j); s_0, s_i)$$
$$+ g(\sim(t(i,j); s_0, s_i), t(i,j); x_0, x_i)\cdot v_{i,t(i,j)},$$
$j = 1, 2, \ldots, m(i)$; $i = 1, 2, \ldots, n$, where $h$ is a known strictly monotone transformation indexed by a fixed-effects parameter vector $\gamma_0$ and a random-effects parameter vector $\gamma_i$, $\gamma_i$'s are *i.i.d.* with mean 0 and unknown covariance matrix $d_\gamma$, $v_{i,t(i,j)}^*$'s are independent errors with mean 0, $\sim$ is a known positive regression function indexed by a fixed-effects parameter vector $s_0$ and a random-effects parameter vector $s_i$, $s_i$'s are *i.i.d.* with mean 0 and unknown covariance matrix $d_s$, $\sim(t; s_0, s_i) = h(x_{i,t}; \gamma_0, \gamma_i)$ for $t \geq 0$, $g$ is a known positive weight function indexed by a fixed-effects parameter vector $x_0$ and a random-effects parameter vector $x_i$, $x_i$'s are *i.i.d.* with mean 0

and unknown covariance matrix $\sigma_{X_i}$ and $V_{i,f(i,j)}$'s are independent standardized errors with mean 0 and variance 1, and independent of $S_i$'s, $\lambda_i$'s and $X_i$'s.

For each fixed $(\lambda_0, \lambda_i)$ pair, let $h^{-1}(\cdot; \lambda_0, \lambda_i)$ be the inverse function of $h(\cdot; \lambda_0, \lambda_i)$ and set $(h^{-1})^{(k)}(v; \lambda_0, \lambda_i) = \partial^k h^{-1}(v; \lambda_0, \lambda_i)/\partial v^k$ for $k = 1, 2, \ldots$. Note that, for each $(i,j)$ pair, both the conditional mean and median of $y_{i,f(i,j)}$ given $S_i$, $\lambda_i$ and $X_i$ are in general unavailable.

As an example,
$$h(u, \lambda_0, \lambda_i) = u^{(\lambda_0 + \lambda_i)}, \; u > 0,$$
$$\sim(t; S_0, S_i) = (S_{01} + S_{i1}) + (S_{02} + S_{i2}) \cdot t, \; t \geq 0,$$
and
$$g(\sim(t; S_0, S_i), t; X_0, X_i)$$
$$= \exp\{X_{01} + X_{02} \cdot \log[\sim(t; S_0, S_i)] + X_i\}, \; t \geq 0,$$
$i = 1, 2, \ldots, n$, can be used for modeling positive longitudinal data, where $S_0 = (S_{01}, S_{02})$ and $X_0 = (X_{01}, X_{02})$. As another example,
$$h(u, \lambda_0, \lambda_i) = u^{[\lambda_0 + \lambda_i]}, \; u \in \boldsymbol{R},$$
$$\sim(t; S_0, S_i) = (S_{01} + S_{i1}) + (S_{02} + S_{i2}) \cdot t, \; t \geq 0,$$
and
$$g(\sim(t; S_0, S_i), t; X_0, X_i)$$
$$= \exp[X_{01} + X_{02} \cdot |\sim(t; S_0, S_i)| + X_i], \; t \geq 0,$$
$i = 1, 2, \ldots, n$, can be used for modeling real-valued longitudinal data, where $S_0 = (S_{01}, S_{02})$ and $X_0 = (X_{01}, X_{02})$.

In the following, we propose an estimation method for all parameters and a time-to-failure distribution in the GTP model with heteroscedastic errors as follows:

To simplify the notation, for each $(i,j)$ pair and $k = 1, 2, \ldots$, set $S = (S_0^T, S_1^T, \ldots, S_n^T)^T$, $\lambda = (\lambda_0^T, \lambda_1^T, \ldots, \lambda_n^T)^T$, $X = (X_0^T, X_1^T, \ldots, X_n^T)^T$, $\theta = (S^T, \lambda^T, X^T)^T$, $\sim_i(S) = \sim(f(i,j); S_0, S_i)$, $h_{ij}(\lambda) = h(y_{i,f(i,j)}; \lambda_0, \lambda_i)$, $g_{ij}(S, X) = g(\sim_i(S), f(i,j); X_0, X_i)$, $H_{ij1}(\theta) = h_{ij}(\lambda) - \sim_i(S, \lambda)$, $H_{ij2}(\theta) = H_{ij1}^2(\theta) - g_{ij}^2(S, X)$, $H_{ij}(\theta) = (H_{ij1}(\theta), H_{ij2}(\theta))^T$, $h^{(k)}(u, \lambda_0, \lambda_i) = \partial^k h(u, \lambda_0, \lambda_i)/\partial u^k$, $g^{(k)}(u, f(i,j); X_0, X_i) = \partial^k g(u, f(i,j); X_0, X_i)/\partial u^k$, $h_\lambda(u, \lambda_0, \lambda_i) = \partial h(u, \lambda_0, \lambda_i)/\partial \lambda$, $h_\lambda^{(k)}(u, \lambda_0, \lambda_i) = \partial^k h_\lambda(u, \lambda_0, \lambda_i)/\partial u^k$, $h_0(v, \lambda_0, \lambda_i) = h_\lambda(h^{-1}(v, \lambda_0, \lambda_i); \lambda_0, \lambda_i)$, $h_0^{(0)}(v, \lambda_0, \lambda_i) = h_0(v, \lambda_0, \lambda_i)$ and $h_0^{(k)}(v, \lambda_0, \lambda_i) = \partial^k h_0(v, \lambda_0, \lambda_i)/\partial v^k$.

By the method of fixed-sample optimal estimating functions (Heyde, 1988), we may wish to utilize the root $\theta^*$ of the fixed-sample optimal estimating equation $S(\theta)|_{\theta = \theta^*} = 0$ to estimate $\theta$, where
$$S(\theta)$$
$$= \Sigma_{i=1}^n \Sigma_{j=1}^{m(i)} E(-\partial H_{ij}^T(\theta)/\partial \theta | S_i, \lambda_i, X_i)$$
$$Cov^{-1}(H_{ij}(\theta) | S_i, \lambda_i, X_i) H_{ij}(\theta).$$

Since both $E(-\partial H_{ij}^T(\theta)/\partial \theta | S_i, \lambda_i, X_i)$ and $Cov(H_{ij}(\theta) | S_i, \lambda_i, X_i)$ are unavailable for each $(i,j)$ pair in the GTP model with heteroscedastic errors, both $S(\theta)$ and $\theta^*$ are unavailable.

By the method of generalized estimating equations (GEE) (Liang and Zeger, 1986), we may wish to utilize the root $\theta^{**}$ of the generalized estimating equation $G^{**}(\theta)|_{\theta = \theta^{**}} = 0$ to estimate $\theta$ instead of the unavailable estimator $\theta^*$, where
$$G^{**}(\theta)$$
$$= \Sigma_{i=1}^n \Sigma_{j=1}^{m(i)} E(-\partial H_{ij}^T(\theta)/\partial \theta | S_i, \lambda_i, X_i) W_{ij}^{-1}(\theta)$$
$$H_{ij}(\theta),$$
with working covariance matrices $W_{ij}(\theta) = \text{diag}\{g_{ij}^2(S, X), 2 \cdot g_{ij}^4(S, X)\}$. Note that if all standardized errors have the same first four moments as those of the standard normal distribution, then $W_{ij}(\theta) = Cov(H_{ij}(\theta) | S_i, \lambda_i, X_i)$ and thus both $\theta^*$ and $\theta^{**}$ are the same. Since $E(-\partial H_{ij}^T(\theta)/\partial \theta | S_i, \lambda_i, X_i)$ is unavailable for each $(i,j)$ pair, both $G^{**}(\theta)$ and $\theta^{**}$ are in general unavailable.

For each $(i,j)$ pair, by deleting all of the terms involving unmodeled $E(V_{i,f(i,j)}^3)$, $E(V_{i,f(i,j)}^4)$, …, we may approximate both $E(\partial h_{ij}(\lambda)/\partial \lambda | S_i, \lambda_i, X_i)$ and $E(H_{ij1}(\theta) \cdot \partial h_{ij}(\lambda)/\partial \lambda | S_i, \lambda_i, X_i)$ by the Taylor approximation, respectively. Then, by the method of unbiased estimating functions (e.g., see Godambe, 1991), we may utilize the root $\theta^\wedge$ of the unbiased estimating equation $G(\theta)|_{\theta = \theta^\wedge} = 0$ to estimate $\theta$ instead of the generally unavailable estimator $\theta^{**}$, where
$$G(\theta) = \Sigma_{i=1}^n \Sigma_{j=1}^{m(i)} D_{ij}^T(\theta) W_{ij}^{-1}(\theta) H_{ij}(\theta),$$
where $D_{ij}(\theta)$ is the Taylor approximation of $E(-\partial H_{ij}^T(\theta)/\partial \theta | S_i, \lambda_i, X_i)$.

One way to obtain the estimator $\theta^\wedge$ is to utilize the following iteration method: We first choose a good initial value $\theta^{\wedge(0)}$ and then iterate the following equations
$$\theta^{\wedge(k+1)} = \theta^{\wedge(k)} + J^1(\theta^{\wedge(k)}) G(\theta^{\wedge(k+1)}),$$
$k = 0, 1, 2, \ldots$, until $\theta^{\wedge(k)}$'s converge to $\theta^\wedge$, where
$$J^1(\theta) = \Sigma_{i=1}^n \Sigma_{j=1}^{m(i)} D_{ij}^T(\theta) W_{ij}^{-1}(\theta) D_{ij}(\theta).$$

Finally, the remaining parameters $d_s$, $d_j$ and $d_x$ could be estimated by $\hat{d_s} = (\sum_{i=1}^{n} \hat{s_i}\hat{s_i}^T)/n$, $\hat{d_j} = (\sum_{i=1}^{n} \hat{j_i}\hat{j_i}^T)/n$, $\hat{d_x} = (\sum_{i=1}^{n} \hat{x_i}\hat{x_i}^T)/n$, respectively.

For each subject $i$, the failure time $T[i]$ is defined as the time when the actual path of subject $i$ crosses the critical level $D$ for some known $D > 0$. Since we only observe the sample paths of subjects, we never observe the actual failure times. For each subject $i$, we may utilize the root $\hat{T[i]}$ of equation $h^{-1}(\sim(t;s_0,\hat{s_i});\hat{j_0},\hat{j_i})|_{t=\hat{T[i]}} = D$ to predict the unobserved failure time $T[i]$.

In some longitudinal studies, we are interested in estimating the time-to-failure distribution $F_D$ for the population of all possible subjects. Since the failure times $T[i]$'s are *i.i.d.* $F_D$, we may utilize the estimator $\hat{F_D}$ to estimate the time-to-failure distribution $F_D$, where

$$\hat{F_D}(t) = [\sum_{i=1}^{n} 1_{(-\infty,t]}(\hat{T[i]})]/n, \ t \in \mathbf{R}.$$

## 四、計畫成果自評

Sometimes, it is inappropriate to assume that the GP model holds for modeling longitudinal data. For example, in practice the transformation $h$ in the GP model is in general unknown. Moreover, different subjects may need different transformations for modeling longitudinal data. Thus, it is better to assume that the transformation $h$ in the GP model depend on fixed-effects and/or random-effects parameters for modeling longitudinal data. For some longitudinal data, we may need not only a transformation but also a weight function in order approximately to achieve homogeneity. This is particularly true when the variance depends on time. Hence, a weight function depending on time $t$ may be needed for transformed longitudinal data. Therefore, it is better to utilize the GTP model with heteroscedastic errors for modeling longitudinal data instead of the GP model.

## 五、參考文獻

[1] Atkinson, A. C. and Cox, D. R. (1988). Transformations (update). In *Encyclopedia of Statistical Sciences*, Volume 9, Kotz, S. and Johnson, N. L. (Ed.), John Wiley & Sons, New York.

[2] Box, G. E. P. and Cox, D. R. (1964). An analysis of transformations. *J. Roy. Statist. Soc. Ser. B* 26, 211-243.

[3] Dowling, N. E., (1993). *Mechanical Behavior of Materials.* Prentice Hall, New York.

[4] Godambe, V. P. (Ed.) (1991). *Estimating Functions.* Oxford, New York.

[5] Heyde, C. C. (1988). Fixed sample and asymptotic optimality for classes of estimating functions, *Contemp. Math.* **80**, 241-247.

[6] Liang, K. Y. and Zeger, S. L. (1986). Longitudinal data analysis using generalized linear models. *Biometrika* **73**, 13-22.

[7] Lu, C. J. and Meeker, W. Q. (1993). Using degradation measures to estimate a time-to-failure distribution, *Technometrics* **35**, 161-174.

[8] Taylor, J. M. G. (1998). Transformations (update). In *Encyclopedia of Statistical Sciences*, Update Volume 2, Kotz, S. (Ed.), John Wiley & Sons, New York.

# 行政院國家科學委員會補助專題研究計畫成果報告
※※※※※※※※※※※※※※※※※※※※※※※※※※
※　　　　　　　　　　　　　　　　　　　　　　※
※　　　　　用 Gamma 過程分析單調長期追蹤資料　　　　※
※　　　　　　　　　　　　　　　　　　　　　　※
※※※※※※※※※※※※※※※※※※※※※※※※※※

計畫類別：□個別型計畫　　□整合型計畫

計畫編號：NSC 90-2118-M-009-009

執行期間：90 年 8 月 1 日至 91 年 7 月 31 日

計畫主持人：陳志榮

共同主持人：

計畫參與人員：

本成果報告包括以下應繳交之附件：
　　□赴國外出差或研習心得報告一份
　　□赴大陸地區出差或研習心得報告一份
　　□出席國際學術會議心得報告及發表之論文各一份
　　□國際合作研究計畫國外研究報告書一份

執行單位：國立交通大學統計所

中　華　民　國　91 年 10　月　30　日