

Application of Local Influence Diagnostics to the Buckley-James Model

Nazrina Aziz¹ and Dong Q Wang²

¹Universiti Utara Malaysia

²Victoria University of Wellington, New Zealand

19th International Conference on Computational Statistics
Paris, France
August 22-28, 2010

Introduction

- Buckley-James Model?

Method used to resolve the problem of a data set containing censored observations

- Censored observations?

A data set that contains observations with incomplete information occurs when the event of interest is not observed

- Some examples:

- In biological research: the time from diagnosis to death
- For industrial research: an example of special interest can be the life time of machine components.

- Are there any other approaches?

Methods used are based on regression ideas, i.e the Miller's method, the Cox method and Koul-Susarla-Van Ryzin estimators

Introduction

Buckley-James censored regression

Multivariate censored regression

Previous Diagnostics Analysis for the BJ method

Local Influence Diagnostics for BJ model

Conclusion

Motivation

Objective of the study

Outline

1

Introduction

• Motivation

- Objective of the study

2

Buckley-James censored regression

3

Multivariate censored regression

4

Previous Diagnostics Analysis for the BJ method

5

Local Influence Diagnostics for BJ model

- Continue...
- Perturbing the variance
- Continue...Perturbing the variance
- Perturbing the response variables
- Perturbing independent variables
- Illustration

• Continue... Illustration

Motivation

- Which methods perform better?

Study	Preference
Miller & Halpern (1982)	The Buckley-James method
Heller & Simonoff (1990)	The Buckley-James method
Heller & Simonoff (1992)	The Buckley-James & Cox method
Stare, Heinzl & Harrell (2000)	The Buckley-James method

- Current study: Buckley-James method
- But it is rarely used. Why?
 - Limited diagnostics analysis developed for the Buckley-James method
 - In the previous diagnostics of Buckley-James model, influential observations merely come from uncensored observations in the data set.

Outline

- 1 Introduction
 - Motivation
 - **Objective of the study**
- 2 Buckley-James censored regression
- 3 Multivariate censored regression
- 4 Previous Diagnostics Analysis for the BJ method
- 5 Local Influence Diagnostics for BJ model
 - Continue...
 - Perturbing the variance
 - Continue...Perturbing the variance
 - Perturbing the response variables
 - Perturbing independent variables
 - Illustration
 - Continue... Illustration

Objective of the study

- **Solution?**

Current study is designed to develop a diagnostic tool for the Buckley-James method

- **What is the new diagnostic tool?**

The local influence diagnostics for the Buckley-James model, which consist of

- variance perturbation;
- response variable perturbation;
- independent variables perturbation.

- **Advantages**

The proposed diagnostics improves the previous ones by taking into account both censored and uncensored data to have a possibility to become an influential observation

Buckley-James censored regression

- Who introduced Buckley-James method?

Buckley and James in 1979

- How does it works?

- Modify the least square standard equations to make it suitable for a data set exposed to censored observations.
- First, review the standard linear regression with the complete data set:

$$Y_i = \alpha + \beta x_i + \varepsilon_i$$

- Let the i th observation have a related censoring time, t_i .
- Now observe Z_i , δ_i and x_i for $i = 1, 2, \dots, n$ where

$$Z_i = \min(y_i, t_i)$$

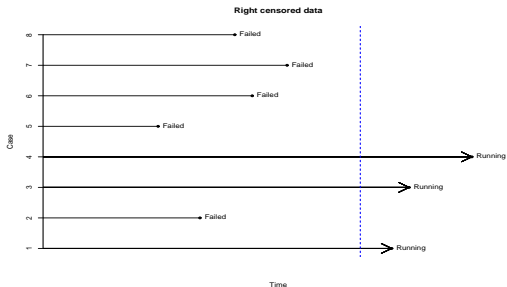


Figure: Plot of right censored data with the dashed lines representing the time line

- and

$$\delta_i = \begin{cases} 0 & \text{(censored) if } y_i \geq t_i, \\ 1 & \text{(uncensored) if } y_i < t_i \end{cases}$$

- Renovate the old response variable (survival time) based on its censored status, δ_j .

$$y_i^*(b) = \begin{cases} bx_i + [\epsilon_i(b)\delta_i + \hat{E}_b(\epsilon_i(b)|\epsilon_i(b) > c_i(b))(1 - \delta_i)] & \text{if } \delta_i = 0, \\ y_i & \text{if } \delta_i = 1 \end{cases}$$

- The residual is represented by the different types of notation $c_i(b) = t_i - bx_i$ or $\epsilon_i(b) = y_i - bx_i$. Choose $e_i(b) = \min\{c_i(b), \epsilon_i(b)\}$
- Note that

$$\begin{aligned} \hat{E}_b(\epsilon_i(b)|\epsilon_i(b) > c_i(b)) &= \frac{\int_{e_i}^{\infty} \epsilon d\hat{F}_b(\epsilon)}{\int_{e_i}^{\infty} d\hat{F}_b(\epsilon)} \\ &= \sum_{k=1}^n w_{ik}(b)e_k(b) \end{aligned}$$

- Next, one can develop the Buckley-James estimator of β as follows

$$\sum_{i=1}^n (x_i - \bar{x})(Y_i^* - x_i\beta) = 0.$$

- By using the iteration, first get the initial estimate of the slope, $\hat{\beta}^{(0)}$, then the Buckley-James estimator of β can be obtained as below

$$\frac{\sum_{i=1}^n (x_i - \bar{x})Y_i^*(\hat{\beta}_n)}{\sum_{i=1}^n (x_i - \bar{x})^2} = \hat{\beta}_{n+1}$$

where $\hat{\beta}_n$ is the estimate of β for the n th iteration, $n = 1, 2, \dots$

- The iteration is stopped when $|\hat{\beta}_{n+1} - \hat{\beta}_n|$ is small and reaches convergence.
- Later one can estimate $\hat{\alpha}$ as follows

$$\hat{\alpha} = \frac{Y^*(\hat{\beta}) - \hat{\beta}x_i}{n} \quad (1)$$

Multivariate censored regression

- What about multivariate censored regression?

Consider $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim F$

- How does it work?

- First, the renovated response variable needs to be obtained as the linear censored regression,

$$\mathbf{Y}^*(\mathbf{b}) = \mathbf{X}\mathbf{b} + \mathbf{W}(\mathbf{b})(\mathbf{Z} - \mathbf{X}\mathbf{b})$$

- Next, the Buckley-James estimators can be developed as follows

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{Y}^*$$

$$\mathbf{W}(\mathbf{b}) = \begin{pmatrix} \delta_1 & w_{12}(\mathbf{b}) & w_{13}(\mathbf{b}) & \dots & w_{1n}(\mathbf{b}) \\ 0 & \delta_2 & w_{23}(\mathbf{b}) & \dots & w_{2n}(\mathbf{b}) \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \ddots & w_{(n-1)n}(\mathbf{b}) \\ 0 & 0 & 0 & \dots & \delta_n \end{pmatrix}$$

- and

$$w_{ik}(\mathbf{b}) = \begin{cases} \frac{d\hat{F}(e_k(\mathbf{b}))\delta_k(1 - \delta_i)}{\hat{S}(e_i(\mathbf{b}))} & \text{if } k > i, \\ 0 & \text{if otherwise} \end{cases}$$

- In multivariate censored regression, the iteration concept is still applied to develop the Buckley-James estimators:

$$b_{n+1} = (X^T X)^{-1} X^T (Xb_n + W(b_n)(Z - Xb_n))$$

- Nevertheless if the iteration fails to converge, one can solve this problem by taking the average of all possible solutions of β

Previous Diagnostics

Renovated Diagnostics	Main purpose
Scatterplot (Smith & Zhang, 1995)	Describe relationship
Hat Matrix (Smith & Zhang, 1995)	Identify outliers
Added Variable Plot (Smith & Peiris, 1999)	Checking linearity Indicates new variable effects
Partial Residual Plot (Wang, Smith & Aziz, 2009)	Independent variable transformations
Cook's Distance (Aziz & Wang, 2009)	Discover influential observations

Local Influence Diagnostics

Local Influence?

Local influence was proposed by Cook in 1986 and it can be used to discover influential observations in a data set.

How does it work?

- The general influence function of $T \in R^{p+1}$, can be displayed as

$$GIF(T, h) = \lim_{\epsilon \rightarrow 0} \frac{T(w_0 + \epsilon h) - T(w_0)}{\epsilon}$$

where $w = w_0 + \epsilon h \in R^n$ describes a perturbation with the null perturbation, w_0 fulfils $T(w_0) = T$ and $h \in R^n$ refers to a unit-length vector.

- Next, one can specify generalised Cook statistics to measure the influence of the perturbations on T as

$$GC(T, h) = \frac{\{GIF(T, h)\}^T M \{GIF(T, h)\}}{c},$$

where M is a $p \times p$ positive-definite matrix and c is a scalar.

Outline

- 1 Introduction
 - Motivation
 - Objective of the study
- 2 Buckley-James censored regression
- 3 Multivariate censored regression
- 4 Previous Diagnostics Analysis for the BJ method
- 5 **Local Influence Diagnostics for BJ model**
 - **Continue...**
 - Perturbing the variance
 - Continue...Perturbing the variance
 - Perturbing the response variables
 - Perturbing independent variables
 - Illustration
 - Continue...Illustration

Continue...

- Continue...How does it work?
 - One may find a direction of $h_{max}(T)$ to perturb a datum and maximize local change in T .
 - The direction of $h_{max}(T)$ can be derived by maximizing the absolute value of $GC(T, h)$ with respect to h .
 - The serious local influence appears if maximum value $GC_{max}(T) = GC(T, h_{max}(T))$.

Outline

- 1 Introduction
 - Motivation
 - Objective of the study
- 2 Buckley-James censored regression
- 3 Multivariate censored regression
- 4 Previous Diagnostics Analysis for the BJ method
- 5 **Local Influence Diagnostics for BJ model**
 - Continue...
 - **Perturbing the variance**
 - Continue...Perturbing the variance
 - Perturbing the response variables
 - Perturbing independent variables
 - Illustration
 - Continue...Illustration

Perturbing the variance

- By using the Buckley-James estimators as follows

$$b = (X^T QX)^{-1} X^T QY^* \quad (2)$$

perturb the variance of the error in (2), by replacing ϵ as $\epsilon_W \sim N(0, \sigma^2 W^{-1})$.

- Let W be diagonal matrix

$$W = \begin{pmatrix} w_1 & 0 & \cdots & 0 \\ 0 & w_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & w_n \end{pmatrix}$$

and vector $w^T = (w_1, w_2, \dots, w_n)$ and w is given by $w = w_o + \epsilon h$, where $w_o^T = (1, 1, \dots, 1)$, the n -vector of ones and $h^T = (h_1, h_2, \dots, h_n)$ refers to a unit-length vector.

- Hence, W can be written as

$$W = I_n + \epsilon D(h), \quad (3)$$

$$\text{where } I_n = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} \text{ and } D(h) = \begin{pmatrix} h_1 & 0 & \cdots & 0 \\ 0 & h_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & h_n \end{pmatrix}.$$

Outline

- 1 Introduction
 - Motivation
 - Objective of the study
- 2 Buckley-James censored regression
- 3 Multivariate censored regression
- 4 Previous Diagnostics Analysis for the BJ method
- 5 Local Influence Diagnostics for BJ model**
 - Continue...
 - Perturbing the variance
 - Continue...Perturbing the variance**
 - Perturbing the response variables
 - Perturbing independent variables
 - Illustration
 - Continue...Illustration

Continue...Perturbing the variance

- Now (2) becomes

$$b(w) = (X^T W Q X)^{-1} X^T W Q Y^*. \quad (4)$$

- By replacing $W = \text{diag}(w_1, w_2, \dots, w_n)$ in (4), $b(w)$ can be rewritten as below $b(w) = [(X^T Q X)^{-1} - \epsilon \{ (X^T Q X)^{-1} X^T Q D(h) X (X^T Q X)^{-1} \}] \times X^T W Q Y^*$, where $X^T W Q Y^* = X^T \{I_n + \epsilon D(h)\} Q Y^* = X^T Q Y^* + \epsilon X^T Q D(h) Y^*$.
- Therefore, $b(w)$ is given by

$$b(w) = b + \epsilon \{ (X^T Q X)^{-1} (X^T Q D(h) e^*) \} + O(\epsilon^2). \quad (5)$$

where $e^* = Y^* - Xb$. From (5), the $GIF(b, h) = (X^T Q X)^{-1} X^T Q D(e^*) h$.

- Next, the generalised Cook statistic of b is developed. It is scaled by $M = X^T \Delta X$, following that $\text{cov}(b) = (X^T \Delta X)^{-1} \sigma_{BJ}^2$, where $\Delta = \text{diag}(\delta_1, \delta_2, \dots, \delta_n)$. Therefore

$$GC_1(b, h) = \frac{h^T D(e^*) (H^*)^2 \Delta D(e^*) h}{ps^2}, \quad (6)$$

where $H^* = X(X^T Q X)^{-1} X^T Q$ and s^2 is the estimate variance.

- By applying $M = X^T X$ to the scaled generalised Cook statistic, which is based on $\text{cov}(b) = (X^T X)^{-1} \sigma^2$,

$$GC_2(b, h) = \frac{h^T D(e^*) (H^*)^2 D(e^*) h}{ps^2}. \quad (7)$$

Outline

- 1 Introduction
 - Motivation
 - Objective of the study
- 2 Buckley-James censored regression
- 3 Multivariate censored regression
- 4 Previous Diagnostics Analysis for the BJ method
- 5 Local Influence Diagnostics for BJ model**
 - Continue...
 - Perturbing the variance
 - Continue...Perturbing the variance
 - Perturbing the response variables**
 - Perturbing independent variables
 - Illustration
 - Continue...Illustration

Perturbing the response variables

- The response variable can be perturbed as $Y_w^* = Y^* + \varepsilon h$, where $h \in R^n$ refers to a unit-length vector. Let equation $(X^T QX)^{-1} X^T QY^*$ become

$$(X^T QX)^{-1} X^T QY_w^* = b + \varepsilon (X^T QX)^{-1} X^T Qh. \quad (8)$$

- Therefore, the general influence function of b under the perturbation can be shown as

$$GIF(b, h) = (X^T QX)^{-1} X^T Qh. \quad (9)$$

- Now two generalised Cook statistics can be developed by using the scale $M = X^T \Delta X$ and $M = X^T X$, which are

$$\text{cov}(b) = \begin{cases} (X^T \Delta X)^{-1} \sigma_{BJ}^2 & \text{if (censored regression),} \\ (X^T X)^{-1} \sigma^2 & \text{if (LSR).} \end{cases} \quad (10)$$

- Hence, $GC_1(b, h) = \frac{h^T (H^*)^2 \Delta h}{ps^2}$ and $GC_2(b, h) = \frac{h^T (H^*)^2 h}{ps^2}$.

Outline

- 1 Introduction
 - Motivation
 - Objective of the study
- 2 Buckley-James censored regression
- 3 Multivariate censored regression
- 4 Previous Diagnostics Analysis for the BJ method
- 5 Local Influence Diagnostics for BJ model**
 - Continue...
 - Perturbing the variance
 - Continue...Perturbing the variance
 - Perturbing the response variables
 - Perturbing independent variables**
 - Illustration
 - Continue...Illustration

Perturbing independent variables

- If one perturbs the i th column of X as $X_W = X + \epsilon l_i h d_i^T$,
- $(X_W^T Q X_W)^{-1} = (X^T Q X)^{-1} - \epsilon l_i (X^T Q X)^{-1} \times (X^T Q h d_i^T + d_i h^T Q X + d_i h^T h d_i^T) (X^T Q X)^{-1} + O(\epsilon^2)$
 and $X_W^T Q Y^* = (X + \epsilon l_i h d_i^T)^T Q Y^* = X^T Q Y^* + \epsilon l_i d_i h^T Q Y^*$.
- Later, $(X_W^T Q X_W)^{-1} (X_W^T Q Y^*) = b + \epsilon l_i (X^T Q X)^{-1} \{d_i h^T Q(e^*) - X^T Q h d_i^T b\} + O(\epsilon^2)$.
- Thus the general influence function of b , $GIF(b, h) = l_i (X^T Q X)^{-1} [d_i h^T Q(e^*) - X^T Q h d_i^T b]$. Replace the i th element of b , therefore $d_i^T b = b_i$ and $GIF(b, h) = l_i (X^T Q X)^{-1} [d_i (e^*)^T - b_i X^T] Q h$.

- Then two generalised Cook statistics for b are constructed as

$$GC_1(b, h) = \frac{l_i^2 h^T H^* \Delta \{e^* d_i^T - b_i X\} (X^T Q X)^{-1} \{d_i (e^*)^T - b_i X^T\} Q h}{ps^2} \text{ and}$$

$$GC_2(b, h) = \frac{l_i^2 h^T H^* \{e^* d_i^T - b_i X\} (X^T Q X)^{-1} \{d_i (e^*)^T - b_i X^T\} Q h}{ps^2}.$$

- One can obtain the diagnostic direction h_{max} by computing the eigenvector corresponding to the largest eigenvalue of the following matrix $H^* \Delta \{e^* d_i^T - b_i X\} (X^T Q X)^{-1} \{d_i (e^*)^T - b_i X^T\} Q$, or $H^* \{e^* d_i^T - b_i X\} (X^T Q X)^{-1} \{d_i (e^*)^T - b_i X^T\} Q$

Outline

- 1 Introduction
 - Motivation
 - Objective of the study
- 2 Buckley-James censored regression
- 3 Multivariate censored regression
- 4 Previous Diagnostics Analysis for the BJ method
- 5 Local Influence Diagnostics for BJ model**
 - Continue...
 - Perturbing the variance
 - Continue...Perturbing the variance
 - Perturbing the response variables
 - Perturbing independent variables
 - Illustration**

Illustration

- Data: Stanford heart transplant data
- $n=152$ patients, 55 deceased ($\delta_i = 1$) while 97 are still alive ($\delta_i = 0$)
- Explanatory variables = censored status, age at time of first transplant (in years) and T5 mismatch score.
- The response variable = survival time(days)
- The Buckley-James model for this data set was developed as $Y = \beta_0 + \beta_1 AGE + \beta_2 AGE^2 + \beta_3 T5$.
- First, consider the variance perturbation. The index plot of $|h_{max}|$ in the next slide shows patients aged below 20 years as the most influential cases.
- This finding agrees well with Reid and Crepeau (1985), and Pettitt and Daud (1989)

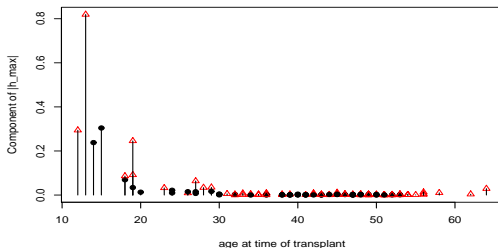


Figure: Index plots of $|h_{max}|$ for perturbing variance for Stanford heart transplant data ($n = 152$).

Outline

- 1 Introduction
 - Motivation
 - Objective of the study
- 2 Buckley-James censored regression
- 3 Multivariate censored regression
- 4 Previous Diagnostics Analysis for the BJ method
- 5 Local Influence Diagnostics for BJ model**
 - Continue...
 - Perturbing the variance
 - Continue...Perturbing the variance
 - Perturbing the response variables
 - Perturbing independent variables
 - Illustration
 - Continue...Illustration

Continue...Illustration

- Next, consider the perturbation of response variable and individual independent variables.
- It is obvious that the most influential patients are aged below 20 years and two patients aged above 60 years.
- Removal of the patients aged 12 and 13 decreases $\hat{\beta}_1$ by 0.010 and 0.030 respectively, while removal of the patient aged 15 increases $\hat{\beta}_1$ by 0.015
- There is no impact on the estimator values in the Buckley-James model when deleting those observations (one at a time) since the maximum eigenvalues for the perturbation of the variance, response variable, x_1 and x_2 are small at 0.142, 0.021, -0.002 and 1.000 respectively.
- However, when the p-value is scrutinized, one can find the p-value for x_1 is roughly five times larger when deleting case 1, and triple when deleting case 4, whereas deleting case 2 has a large effect on the p-value of x_2 where the value becomes fourteen times larger.
- No attention is given to x_3 since this variable is not strongly associated with survival time.

Conclusion

- The proposed local influence diagnostics for the Buckley-James model performs very well for identifying influential cases and for assessing the effects that perturbations to the assumed data would have on inferences.
- It should also be noted that the proposed diagnostics is able to easily detects influential observations from both groups i.e. censored and uncensored observations in the data set as opposed to the previous diagnostics for Buckley-James model.

- Introduction
- Buckley-James censored regression
- Multivariate censored regression
- Previous Diagnostics Analysis for the BJ method
- Local Influence Diagnostics for BJ model
- Conclusion**

Thank you