



Analysis of Influential Observation for Ridge Estimator in Semiparametric Regression

Semra TÜRKAN

Hacettepe University, Department of Statistics, Beytepe, 06800Ankara, Turkey

ABSTRACT

The detection of influential observations has attracted a great deal of attention in last few decades. Most of the ideas of determining influential observations are based on single-case diagnostics with i th case deleted. The Cook's distance is most commonly used among the other single-case diagnostics and successfully applied to various statistical models. In this article, we propose Cook's distance for the ridge regression estimator of the parametric component in the semiparametric regression model to detect influential observations. We investigate the performance of proposed diagnostic to detect influential observations by using simulation data.

KEYWORDS : Semiparametric regression model, Ridge regression estimator, Cook's distance, Influential observations.

1. Introduction

Regression diagnostics consist of a collection of methods used in the identification of influential points and multicollinearity. Particularly, the detection of influential observations has received a great deal of attention in the statistical literature. Hence a number of statistical measures have been proposed to identify influential observations, in which Cook's distance proposed by Cook (1977) is one of the most commonly used influence measures. However, most of the regression diagnostics have been related to parametric regression models. Diagnostic measures in various nonparametric regression or semiparametric regression models are quite rare (Zhang et al., 2007).

In this paper, we examine the influence of observations on the ridge estimator for the vector of parameters β in a semiparametric regression model. Consider the following semiparametric regression model:

$$y_i = \mathbf{x}_i^T \beta + f(t_i) + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (1)$$

where the y_i 's are observations, \mathbf{x}_i is $p \times 1$ vector ($p \leq n$), t_i is scalar, $\beta = (\beta_1, \dots, \beta_p)^T$ is a vector of unknown parameters, $f(\cdot)$ is unknown function, and ε_i 's are independent and identically distributed random variables with zero mean and variance σ^2 .

Semiparametric models are more flexible than linear models since they consist of a parametric and a nonparametric component. These models are used when the response y linearly depends on x , but it is nonlinearly related to t . Hence $f(t)$ represents a smooth unparametrized functional relationship. The aim is to estimate β and nonparametric function $f(t)$ from the data $\{y_i, x_i, t_i\}$.

The model in (1) can be written in matrix-vector notation as

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{f} + \boldsymbol{\varepsilon} \quad (2)$$

where $\mathbf{y} = (y_1, \dots, y_n)^T$, $\mathbf{f} = (f(t_1), \dots, f(t_n))^T$, $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)$ and $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ is the $n \times p$

matrix without any restrained conditions, namely, $rank(X) < p$ or $rank(X) = p$ (ill-conditioned or not).

In fact, if X is an ill-conditioned matrix, then the results may not reliable. Many studies do not consider the case $rank(X) < p$, and the situation that the design matrix X is rank-deficient is rarely investigated. It is noticeable that ridge estimation not only solves rank-deficient and ill-conditioned problems, but also presents a new method which can deal with (non)linear and semiparametric regression models for $rank(X) = p$ without ill-conditioning (Hu 2005). Hence, Hu (2005) proposed ridge estimator of semiparametric regression model. Roozbeh et al. (2010) introduced a semiparametric ridge regression estimator for the vector-parameter when the matrix $X^T X$ is ill-conditioned. The large condition number of parametric component indicates that a ridge regression estimator can be used for β . The reduction in multicollinearity should be a first step for the effective detection of influential observations.

In this article, the Cook's distance is defined to detect influential observations on the ridge estimator of β in semiparametric regression. Approximate deletion formula as functions of corresponding residuals and leverages are proposed. In Section 2, the ridge regression estimator is given. In Section 3, the approximate deletion formula of Cook's distance for the ridge estimator of β in semiparametric regression is derived. In section 4, the simulation data is examined using proposed diagnostic method.

The Model and Ridge Estimator

In this article, the partial kernel smoothing estimator of β , which attains the usual parametric convergence rate $n^{1/2}$ without undersmoothing the nonparametric component $f(\cdot)$, is used. If β is known, a naturel nonparametric estimator of $f(\cdot)$ is

$$\hat{f}(t, \beta) = \sum_{i=1}^n W_{ni}(t)(y_i - x_i^T \beta) \tag{3}$$

where $W_{ni}(\cdot)$ is the positive weight function (Roozbeh et al., 2010).

To estimate β in (2) using kernel weight functions, least square estimator can be used. Consider the following objective function:

$$SS(\beta) = (\tilde{y} - \tilde{X}\beta)^T (\tilde{y} - \tilde{X}\beta) \tag{4}$$

where $\tilde{y} = (\tilde{y}_1, \dots, \tilde{y}_n)$, $\tilde{X} = (x_1^T, \dots, x_n^T)$, $\tilde{y}_i = y_i - \sum_{j=1}^n W_{nj}(t_i)y_j$ and $\tilde{x}_i = x_i - \sum_{j=1}^n W_{nj}(t_i)x_j$ for $i=1, \dots, n$. The first-order condition of objective function (4) minimized by the vector β is obtained as:

$$\hat{\beta} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{y} \tag{5}$$

The properties of the least square estimation of β depend heavily on the characteristics of the $\tilde{X}^T \tilde{X}$. If the $\tilde{X}^T \tilde{X}$ matrix is ill-conditioned (close dependency among various columns of $\tilde{X}^T \tilde{X}$), the least square estimation produce unduly large sampling variances. As a remedy, Hu

(2005) suggested to use ridge regression estimator instead of least square estimator. The ridge regression estimator of β can be obtained by using following objective function:

$$SS(\beta) = (\tilde{y} - \tilde{X}\beta)^T (\tilde{y} - \tilde{X}\beta) + k\beta^T \beta \tag{6}$$

The first-order condition of objective function (6) minimized by the vector β is obtained as:

$$\hat{\beta}_R = (\tilde{X}^T \tilde{X} + kI)^{-1} \tilde{X}^T \tilde{y} \tag{7}$$

where k is a biasing parameter ($k > 0$) and $\tilde{X} = (I - W)X$, $\tilde{y} = (I - W)y$. For the model in (2), the estimates of β and f can be written as

$$\hat{\beta}_R = Ay \tag{8}$$

$$\hat{f}_R = W(y - X\hat{\beta}_R) \tag{9}$$

where $A = [X^T(I - W)^T(I - W)X + kI]^{-1} X^T(I - W)^T(I - W)$ and W is kernel weights matrix (Tabakan and Akdeniz, 2009).

Cook's Distance for Ridge Estimator

In this section, we define Cook's distance for ridge estimator of β in semiparametric regression to gauge the influential observations and express it as function of the corresponding residuals and leverages.

3.1. Influence on $\hat{\beta}$

Following the study of Walker and Birch (1988), an influence measure for the i th observation on $\hat{\beta}$ can be defined as a type of Cook's distance for ridge regression

$$\tilde{C}_i = \frac{(\hat{\beta} - \hat{\beta}_{R,-i})^T \tilde{X}^T \tilde{X} (\hat{\beta} - \hat{\beta}_{R,-i})}{\sigma^2 \text{tr}(\tilde{H})} \tag{10}$$

where $\tilde{H} = \tilde{X}(\tilde{X}^T \tilde{X} + kI)^{-1} \tilde{X}^T$ is matrix which plays the same role as the hat matrix in least square regression. Noting $\text{tr}(\tilde{H}) = p$, the above defined \tilde{C}_i can be expressed as a function of the corresponding residuals and leverages,

$$\tilde{C}_i = \frac{\tilde{e}_i^2}{p\sigma^2} \left[\frac{\sum_{i=1}^n \tilde{h}_{ij}^2}{(1 - \tilde{h}_i)^2} \right] \tag{11}$$

where \tilde{e}_i is the i th component of the residual vector $\tilde{e} = \tilde{y} - \hat{y}_R = \tilde{y} - \tilde{X}\hat{\beta}_R$, \tilde{h}_i is i th diagonal component of \tilde{H} and \tilde{h}_{ij} is the ij th element of \tilde{H} . To compute \tilde{C}_i , it is required an estimator

of σ^2 . A simple estimator is $s^2 = \sum_{i=1}^n \tilde{e}_i^2 / (n - p)$.

3.2. Influence on \hat{y}

An influence measure for the i th observation on the vector of fitted values can be similarly defined by

$$C_i = \frac{(\hat{y}_i - \hat{y}_{R,-i})}{SE(\hat{y}_i)} = \frac{\tilde{\mathbf{x}}_i^T (\hat{\boldsymbol{\beta}}_R - \hat{\boldsymbol{\beta}}_{R,-i})}{SE(\hat{y}_i)} \tag{12}$$

where $SE(\hat{y}_i) = s [\mathbf{x}_i^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + k\mathbf{I})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + k\mathbf{I})^{-1} \mathbf{x}_i]^1/2$ is an estimator of standard error of the fitted value.

C_i can be expressed as a function of the corresponding residuals and leverages as

$$C_i = \frac{1}{s (\sum_{i=1}^n \tilde{h}_{ij}^2)^1/2} \frac{\tilde{h}_i \tilde{e}_i}{(1 - \tilde{h}_i)} \tag{13}$$

The main advantage of deletion formulas is that, as in LS, the estimator does not have to be recomputed every time a case is deleted. For a value of k and h , all of the elements of (11) and (13) are readily available from a single run of semiparametric regression.

Simulation Study

In this section, we consider the model generating the data set which is used in the study of Roozbeh et al. (2010). However, we generate the data that will be included influential observations. The following model is used to generate data:

$$y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + x_{i3}\beta_3 + x_{i4}\beta_4 + x_{i5}\beta_5 + f(t_i) + \varepsilon_i, \quad i = 1, \dots, n$$

where $\boldsymbol{\beta} = (1, 2, 2, -5, 4)$, $\sigma = 0.1$ and $f(t_i) = \sqrt{t_i(1-t_i)} \sin(\frac{2.1\pi}{t_i + 0.05})$ for $t_i = (i - 0.5)/n$, $i=1, \dots, n$. The first 450 observations are generated from $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ and $\mathbf{x}_i \sim N_5(\mathbf{0}, \sum_x)$ with

$$\sum_x = \begin{pmatrix} 0.81 & 0.4 & 0.3 & -0.2 & -0.1 \\ 0.4 & 2.25 & 0.4 & 0.3 & -0.2 \\ 0.3 & 0.4 & 1 & 0.4 & 0.3 \\ -0.2 & 0.3 & 0.4 & 0.64 & 0.4 \\ -0.1 & -0.2 & 0.3 & 0.4 & 0.49 \end{pmatrix}$$

and the last 50 observations which are considered as influential are generated from $\epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ where $\sigma = 2$ and $\mathbf{x}_i \sim N_5(\boldsymbol{\mu}, \Sigma_x)$ with $\boldsymbol{\mu} = (10,10,10,10,10)$ and Σ_x .

The weight function $W_{ni}(t_j)$ is defined as

$$W_{ni}(t_j) = \frac{1}{nh_n} K\left(\frac{t_i - t_j}{h_n}\right) = \frac{1}{nh_n} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(t_i - t_j)^2}{2h_n^2}\right\}$$

which is Priestley and Chao’s weight with the Gaussian kernel (Roozbeh et al., 2010).

We use the cross-validation (CV) method to select the optimal bandwidth h_n and biasing parameter k simultaneously, which minimizes the following CV function:

$$CV(h_n) = \frac{1}{n} \sum_{i=1}^n (\tilde{y}_i - \tilde{y}_{-i})^2 = \frac{1}{n} \sum_{i=1}^n (\tilde{e}_i / 1 - \tilde{h}_{ii})^2,$$

Different combinations of k (0.2, 0.4, 0.6, 0.8, 1) and $h_n = (0.01, 0.06, 0.11, 0.16, 0.21, 0.26, \dots, 0.96)$ are used to find k and h_n which minimize the CV simultaneously. The minimum of CV occurred at $k=0.2$ and $h_n = 0.06$.

The ratio of largest eigenvalue to the smallest eigenvalue of matrix $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ is $\lambda_5 / \lambda_1 = 29571.84 / 38.12 = 775.75$ which implies that $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ matrix is ill-contioned or the existence of multicollinearity in the data set. Hence, in this situation it is supported that ridge estimator is used instead of least square estimates.

In order to examine the efficiency of Cook’s distance \tilde{C}_i defined in (10) and C_i defined in (12) in detecting single observation that has large impact on ridge estimator of $\boldsymbol{\beta}$ in semiparametric regression, as stated above, some observations are generated to be influential observations one by one to see whether the defined measures detect it out. For each observation, these measures are calculated and the index plots of these measures are obtained as in Figure 1 and Figure 2.

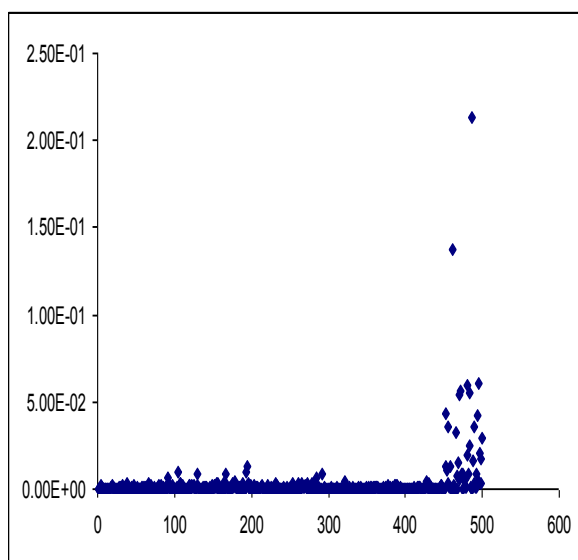


Figure 1. Index plot of \tilde{C}_i

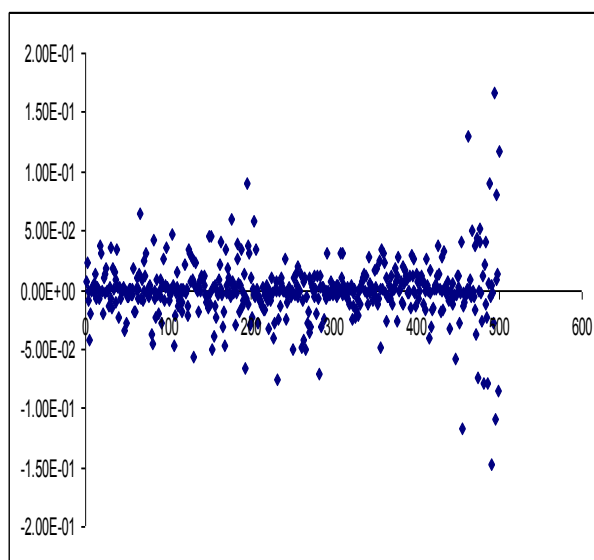


Figure 2. Index plot of C_i

As seen from Figure 1 and Figure 2, \tilde{C}_i can recognize majority of generated the last 50 influential observations while C_i can not as successful as \tilde{C}_i to recognize generated influential observations. In addition, C_i could find some other observations as influential observations.

Conclusion

In this study, Cook's distances that are function of leverages and residuals is studied for ridge estimator in semiparametric regression. Although no conventional cut-off points are introduced or developed for the these measures, it seems that index plot is an optimistic and conventional procedure to disclose influential cases. It is seen that the proposed Cook's distances are successful to detect influential observations in the data.

REFERENCES

- Akdeniz, F., Tabakan, G.: Restricted Ridge Estimators of the Parameters in Semiparametric Regression Model, Communications in Statistics-Theory and Methods, 38: 1852-1869 (2009). Cook, R.D.: Detection of Influential Observations in Linear Regression, Technometrics, 19, 15-18 (1977). Hu, H.: Ridge estimation of semiparametric regression model. J. Comput. Appl. Math. 176, 215-222 (2005). Roozbeh, M., Arashi, M., Niroumand, H.A.: Semiparametric ridge regression approach in partially linear models. Comm. Statist. Simulation Comput. 39, 449-460 (2010). Zhang, C., Mei, C., Zhang, J.: Influence Diagnostics in Partially Varying Coefficient Models, Acta. Math. Appl. Sinica, 23 (4): 619-628 (2007). Walker, E. and Birch, J.B. Influence measures in ridge regression. Technometrics 30, 221-227 (1988).