

Local Empirical Likelihood Diagnosis of Varying Coefficient Density-Ratio Models Based on Case-Control Data

Shuling Wang¹, Lin Zheng², Jiangtao Dai³

¹Department of Fundamental Course, Air Force Logistics College, Xuzhou, China

²School of Statistics and Applied Mathematics, Anhui University of Finance and Economics, Bengbu, China

³Fundamental Science Department, North China Institute of Astronautic Engineering, Langfang, China

Email: 155328313@qq.com

Received 18 August 2014; revised 20 September 2014; accepted 30 September 2014

Copyright © 2014 by authors and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

In this paper, a varying-coefficient density-ratio model for case-control studies is developed. We investigate the local empirical likelihood diagnosis of varying coefficient density-ratio model for case-control data. The local empirical log-likelihood ratios for the nonparametric coefficient functions are introduced. First, the estimation equations based on empirical likelihood method are established. Then, a few of diagnostic statistics are proposed. At last, we also examine the performance of proposed method for finite sample sizes through simulation studies.

Keywords

Varying-Coefficient Density-Ratio Model, Local Empirical Likelihood, Outliers, Influence Analysis

1. Introduction

Varying coefficient models are often used as extensions of classical linear models (e.g. Shumway [1]). Their appeals are that the modeling bias can be significantly reduced and the “curse of dimensionality” can also be avoided. These models have gained considerable attention due to their various applications in many areas, such as biomedical study, finance, econometrics, and environmental study. The estimation for the coefficient functions has been extensively discussed in the literatures, including the smoothing spline method (see Hastie and Tibshirani [2]), the locally weighted polynomial method (see Hoover *et al.* [3]), the two-step estimation procedure (see Fan and Zhang [4]), and the basis function approximations (see Huang *et al.* [5]).

In this paper, we consider the following general two-sample varying-coefficient density-ratio model

$$f(w, z) = \psi \left\{ \alpha(w) + \beta(w)^T z \right\} g(w, z) \tag{1}$$

where $\psi\{\cdot\}$ is a nonnegative known function that makes $f(w, z)$ to be a density function, which includes the exponential-tilt model as a special case with $\psi(\cdot) = \exp(\cdot)$. In parametric situation, Thomas [6] and Lustbader *et al.* [7] considered a general relative risk model, a mixture model, $f(z)/g(z) = (1 + \beta^T z)^\delta \left[\exp\{\beta^T z\} \right]^{1-\delta}$, where δ is a scalar parameter that describes the general shape of the relative risk function. It includes additive relative risk model ($\delta = 1$) and log-linear relative model ($\delta = 0$) as special cases.

Various density-ratio models for some conventional density functions were discussed in Kay and Little [8]. It has been shown recently that the density-ratio model provides a good fit to the observed data in some medical applications (Qin and Zhang [9]; Qin *et al.* [10]; Zhang [11]), genetic quantitative trait loci analysis (Zou *et al.* [12]), and clinical trials with skewed outcomes (White and Thompson [13]). Liu, Jiang and Zhou [14] considered estimation and inference for the two-sample varying-coefficient density-ratio model (1) by constructing the local empirical likelihood function. The EL approach is appealing for analyzing the varying-coefficient density-ratio model because the two density functions in (1) can be modeled nonparametrically. This nonparametric method of inference has sampling properties similar to the bootstrap. Another advantage of the EL approach is that it takes auxiliary information, such as the density-ratio in (1), into account to improve estimation.

The empirical likelihood method origins from Thomas & Grunkemeier [15]. Owen [16] first proposed the definition of empirical likelihood and expounded the system info of empirical likelihood. Zhu and Ibrahim [17] utilized this method for statistical diagnostic. Liugen Xue and Lixing Zhu [18] summarized the application of this method.

Over the last several decades, the diagnosis and influence analysis of linear regression model has been fully developed (R.D. Cook and S. Weisberg [19], Bocheng Wei, Gobin Lu & Jianqing Shi [20]). Regarding the varying coefficient model, especially for the B-spline estimation of parameter, diagnosis and influence analysis have some results (Z. Cai, J. Fan, R. Li [21], J. Fan, W. Zhang [22]). So far the statistical diagnostics of varying-coefficient density-ratio models with case-control data based on local empirical likelihood method has not yet seen in the literature. This paper attempts to study it.

The remainder of the article is organized as follows. Local empirical likelihood and estimation equation are presented in Section 2. The main results are given in Section 3. An example is given to illustrate our results in Section 4.

2. Local Empirical Likelihood and Estimation Equation

Let X_1, \dots, X_{n_1} be a sequence of independent and identically distributed random vectors from the control group, each with density $g(x)$, and $X_{n_1+1}, \dots, X_{n_1+n_2}$ be a sequence of independent and identically distributed random vectors from the case group, each with density $f(x)$, n_1 and n_2 are the number of subjects in the control group and case group, respectively. Let $n = n_1 + n_2$, and $\{X_1, \dots, X_n\} = \{X_1, \dots, X_{n_1}, X_{n_1+1}, \dots, X_{n_1+n_2}\}$ denote the pooled sample. Assume that $n_1/n \rightarrow \rho > 0$ as $n \rightarrow \infty$. From model (1), the empirical likelihood function derived according to Prentice and Pyke [23] is:

$$l(\alpha, \beta, G) = \prod_{j=1}^n dG(X_j) \prod_{i=n_1+1}^n \psi(X_j) dG(X_j) = \prod_{j=1}^n \tilde{p}_j \prod_{i=n_1+1}^n \psi(X_i) \tag{2}$$

where $\psi(x) = \psi(w, z) = \psi(\varphi(x)) = \psi\{\alpha(w) + \beta(w)^T z\}$, $\tilde{p}_i = dG(X_i)$ and $G(X)$ is the distribution function corresponding to $g(x)$. However, $l(\alpha, \beta, G)$ can not be used directly to obtain estimates for $\alpha(\cdot)$ and $\beta(\cdot)$ because $\alpha(\cdot)$ and $\beta(\cdot)$ are infinite-dimensional parameters. Thus, instead of (2), we consider the localized conditional empirical likelihood below.

Assume that all components of $\alpha(\cdot)$ and $\beta(\cdot)$ are smooth so that they admit Taylors series expansions, *i.e.*, for each given w_0 and for w around w_0 ,

$$\begin{aligned} \beta(w) &\approx \beta(w_0) + \beta'(w_0)(w - w_0), \\ \alpha(w) &\approx \alpha(w_0) + \alpha'(w_0)(w - w_0). \end{aligned} \tag{3}$$

Let $\xi(w) = (\alpha(w), \beta(w)^T, \alpha'(w), \beta'(w)^T)^T$, and $X_i^*(w) = (1, Z_i^T, W_i - w, Z_i^T(W_i - w)^T)^T$. For simplicity, denote $\xi(w_0)$ by ξ and $X_i^*(w_0)$ by X_i^* for fixed w_0 . Then, the local log empirical likelihood (LEL) function $l(\xi)$ of is

$$l(\xi) = \sup \left\{ \sum_{i=1}^n w_h(W_i, w_0) \log p_i + \sum_{i=n_1+1}^n w_h(W_i, w_0) \log (\kappa(W_i) \psi(\xi^T X_i^*)) : p_i \geq 0, 1 \leq i \leq n, \sum_{i=1}^n p_i = 1, \sum_{i=1}^n p_i (\kappa(W_i) \psi(\xi^T X_i^*) - 1) = 0 \right\},$$

where $w_h(t_i, t_0) = K_h(t_i - t_0) / \sum_{j=1}^n K_h(t_j - t_0)$ is the weight with kernel function $K_h(\cdot) = K(\cdot/h)/h$ and h represents the size of the local neighborhood. The kernel weight is used to give smoother weight to data with w near w_0 . The last constraint is the auxiliary information for the EL estimation. By the method of Lagrange multipliers, similar to that used in Owen (2001), we obtain

$$p_i = \frac{w_h(W_j, w_0)}{1 + \lambda^T (\kappa(W_i) \psi(\xi^T X_i^*) - 1)}$$

where λ is determined by the constraint equation

$$\sum_{j=1}^n \frac{\kappa(W_i) \psi(\xi^T X_i^*) - 1}{1 + \lambda^T (\kappa(W_i) \psi(\xi^T X_i^*) - 1)} w_h(W_j, w_0) = 0.$$

Motivated by Zhu and Ibrahim (2008), we regard λ and ξ as independent variables and define

$$Q_n(\lambda, \xi) = -n^{-1} \sum_{i=1}^n \log \left(1 + \lambda^T (\kappa(W_i) \phi(\xi^T X_i^*) - 1) \right).$$

Obviously, the maximum empirical likelihood estimates $\hat{\xi}$ and $\hat{\lambda}$ are the solutions of following equations:

$$\begin{cases} Q_{1,n}(\lambda, \xi) = \frac{\partial Q_n(\lambda, \xi)}{\partial \lambda} = -n^{-1} \sum_{i=1}^n (\kappa(W_i) \phi(\xi^T X_i^*) - 1) \left\{ 1 + \lambda^T (\kappa(W_i) \phi(\xi^T X_i^*) - 1) \right\}^{-1} = 0, \\ Q_{2,n}(\lambda, \xi) = \frac{\partial Q_n(\lambda, \xi)}{\partial \xi} = -n^{-1} \sum_{i=1}^n \frac{\partial (\kappa(W_i) \phi(\xi^T X_i^*) - 1)}{\partial \xi} \lambda \left\{ 1 + \lambda^T (\kappa(W_i) \phi(\xi^T X_i^*) - 1) \right\}^{-1} = 0. \end{cases}$$

3. Local Influence Analysis of Model

We consider the local influence method for a case-weight perturbation $\omega \in R^n$, for which the empirical log-likelihood function $l_E(\xi|\omega)$ is defined by $l_E(\xi|\omega) = \sum_{i=1}^n \omega_i l_{E,i}(\xi)$. In this case, $\omega = \omega^0$, defined to be an $n \times 1$ vector with all elements equal to 1, represents no perturbation to the empirical likelihood, because $l_E(\xi|\omega^0) = l_E(\xi)$. Thus, the empirical likelihood displacement is defined as $LD_E(\omega) = 2 \left[l_E(\hat{\xi}|\omega) - l_E(\hat{\xi}|\omega^0) \right]$, where $\hat{\xi}(\omega)$ is the maximum empirical likelihood estimator of ξ based on $l_E(\xi|\omega)$. Let $\omega(a) = \omega^0 + ah$ with $\omega(0) = \omega^0$ and $d\omega(a)/da|_{a=0} = h$, where h is a direction in R^n . Thus, the normal curvature of the influence graph $(\omega^T, LD_E(\omega))^T$ is given by $C_h(\omega^0) = h^T H_{LD_E(\omega^0)} h$, where

$$H_{LD_E(\omega^0)} = -2 \frac{\partial^2 LD_E \{ \hat{\xi}(\omega) \}}{\partial \omega \partial \omega^T} \Big|_{\omega^0} = 2 \Delta^T \left\{ -\partial_{\xi}^2 l_E(\xi) \right\}^{-1} \Delta \Big|_{\omega^0, \hat{\xi}},$$

in which $\Delta = \partial_{\xi \omega}^2 LD_E(\xi, \omega)$ is a $p \times n$ matrix with (k, i) -th element given by $\partial_{\xi_k} l_{E,i}(\xi)$.

We consider two local influence measures based on the normal curvature $C_h(\omega^0)$ as follows. Let $\lambda_1 \geq \dots \geq \lambda_p \geq \lambda_{p+1} = \dots = \lambda_n = 0$ be the ordered eigenvalues of the matrix $H_{LD_E(\omega^0)}$ and let

$\{v_m = (v_{m1}, \dots, v_{mn})^T : m = 1, \dots, n\}$ be the associated orthonormal basis, that is, $H_{LD_E(\omega^0)}v_m = \lambda_m v_m$. Thus, the spectral decomposition of $H_{LD_E(\omega^0)}$ is given by

$$H_{LD_E(\omega^0)} = \sum_{m=1}^n \lambda_m v_m v_m^T .$$

The most popular local influence measures include v_1 , which corresponds the largest eigenvalue λ_1 , as well as $C_{e_j} = \sum_{m=1}^p \lambda_m v_{mj}^2$, where e_j is an $n \times 1$ vector with j -th component 1 and 0 otherwise. The v_1 represents the most influential perturbation to the empirical likelihood function, whereas the j -th observation with a large C_{e_j} can be regarded as influential.

As the discuss of Zhu *et al.* (2008), for varying-coefficient density-ratio model, we can deduce that

$$C_{e_j} = 2ELD_j \{1 + o_p(1)\} = 2ECD_j \{1 + o_p(1)\} = -2n^{-1} \Delta_j^T S_{22.1}^{-1} \Delta_j \{1 + o_p(1)\}, \tag{4}$$

where $\Delta_j = \partial_{\xi} l_{E,j}(\xi) \Big|_{\xi=\hat{\xi}} = \frac{S_{21} S_{11}^{-1} (\kappa(W_j) \psi(\hat{\xi}^T X_j^*) - 1)}{1 + \hat{\lambda}^T (\kappa(W_j) \psi(\hat{\xi}^T X_j^*) - 1)} + o_p(1)$,

$$S_{11} = \partial_{\lambda} Q_{1,n} = \frac{1}{n} \sum_{i=1}^n \frac{(\kappa(W_i) \psi(\xi^T X_i^*) - 1)(\kappa(W_i) \psi(\xi^T X_i^*) - 1)^T}{(1 + \lambda^T (\kappa(W_i) \psi(\xi^T X_i^*) - 1))^2} \Bigg|_{\xi=\hat{\xi}, \lambda=\hat{\lambda}},$$

$$S_{12} = \partial_{\xi} Q_{1,n} = \frac{1}{n} \sum_{i=1}^n \frac{(\kappa(W_i) \psi(\xi^T X_i^*) - 1) \lambda^T \partial_{\xi} (\kappa(W_i) \psi(\xi^T X_i^*) - 1) - \partial_{\xi} (\kappa(W_i) \psi(\xi^T X_i^*) - 1) (1 + \lambda^T (\kappa(W_i) \psi(\xi^T X_i^*) - 1))}{(1 + \lambda^T (\kappa(W_i) \psi(\xi^T X_i^*) - 1))^2} \Bigg|_{\xi=\hat{\xi}, \lambda=\hat{\lambda}},$$

$$S_{21} = \partial_{\lambda} Q_{2,n} = \frac{1}{n} \sum_{i=1}^n \frac{(\kappa(W_i) \psi(\xi^T X_i^*) - 1) \lambda^T \partial_{\xi} (\kappa(W_i) \psi(\xi^T X_i^*) - 1) - \partial_{\xi} (\kappa(W_i) \psi(\xi^T X_i^*) - 1) (1 + \lambda^T (\kappa(W_i) \psi(\xi^T X_i^*) - 1))}{(1 + \lambda^T (\kappa(W_i) \psi(\xi^T X_i^*) - 1))^2} \Bigg|_{\xi=\hat{\xi}, \lambda=\hat{\lambda}},$$

$$S_{22} = \partial_{\xi} Q_{2,n} = \frac{1}{n} \sum_{i=1}^n \frac{\partial_{\xi}^T (\kappa(W_i) \psi(\xi^T X_i^*) - 1) \lambda \lambda^T \partial_{\xi} (\kappa(W_i) \psi(\xi^T X_i^*) - 1)}{(1 + \lambda^T (\kappa(W_i) \psi(\xi^T X_i^*) - 1))^2} \Bigg|_{\xi=\hat{\xi}, \lambda=\hat{\lambda}},$$

$$S_{22.1} = -S_{21} S_{11}^{-1} S_{12}.$$

4. Numerical Study

We generate X_1, \dots, X_{n_1} and $X_{n_1+1}, \dots, X_{n_1+n_2}$ from two densities $f(x)$ and $g(x)$, respectively. We set both densities $f(x)$ and $g(x)$ to be trivariate normal distributions, in which $x = (w, z^T)^T$, W is a scalar, $z^T = (z_1, z_2)$, and

$$g(w, z) = (2\pi)^{-3/2} |\Sigma_g|^{-1/2} \exp \left\{ -2^{-1} \left[(w-1, z_1, z_2) (\Sigma_g)^{-1} (w-1, z_1, z_2)^T \right] \right\},$$

$$f(w, z) = (2\pi)^{-3/2} |\Sigma_f|^{-1/2} \exp \left\{ -2^{-1} \left[(w-1, z_1, z_2 - 1) (\Sigma_f)^{-1} (w-1, z_1, z_2 - 1)^T \right] \right\},$$

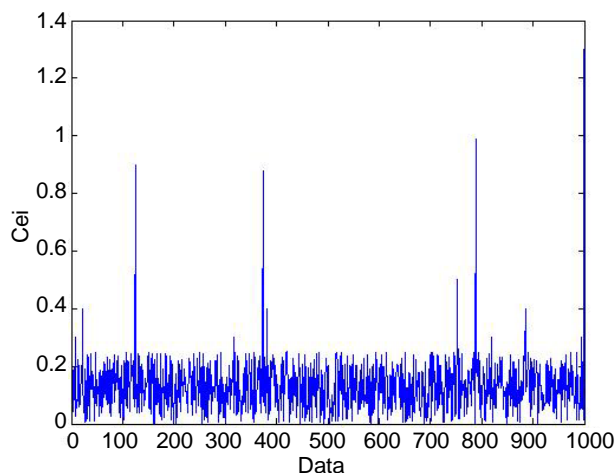


Figure 1. The influence value of C_{e_i} .

are trivariate normal densities with means $\mu_g = (1, 0, 0)^T$ and $\mu_f = (1, 0, 1)^T$, and inverses of the covariances

$$\Sigma_g^{-1} = \begin{pmatrix} 1/2 & -2/3 & 1/3 \\ -2/3 & 2 & 0 \\ 1/3 & 0 & 1 \end{pmatrix}, \quad \Sigma_f^{-1} = \begin{pmatrix} 2/3 & -1/3 & -2/3 \\ -1/3 & 2 & 0 \\ -2/3 & 0 & 1 \end{pmatrix}.$$

Because $f(w, z)/g(w, z) = \exp\{-z_1(w-1)/3 + z_2w - (w-1)^2/12 - 2(w-1)/3 - 1/2\}$, we have $\alpha(w) = -(w-1)^2/12 - 2(w-1)/3 - 1/2$, $\beta_1(w) = -(w-1)/3$ and $\beta_2(w) = w$.

We draw 1000 data sets with sample size $n = n_1 + n_2 = 400$ for various values of $\rho_n = n_1/n = 0.6$. We choose the Epanechnikov kernel $K(t) = 0.75(1-t^2)_+$ to localize the coefficient functions.

In order to check out the validity of our proposed methodology, we change the value of the first, 125th, 374th, 789th and 999th data. For every case, it is easy to obtain $\kappa(W_i)\psi(\xi^T X_i^*) - 1$. For ξ and λ , using the samples, we evaluated their maximum empirical likelihood estimators.

Consequently, it is easy to calculate the value of $S_{11}, S_{12}, S_{21}, S_{22}$ and C_{e_i} . The result of C_{e_i} is as following **Figure 1**.

It can be seen from the result of C_{e_i} that the first, 125th, 374th, 789th and 999th data are strong influence points. Indeed, our results are illustrated.

5. Discussion

In this paper, we considered the statistical diagnosis for varying-coefficient density-ratio model based on local empirical likelihood. Through simulation study, we illustrate that our proposed method can work fairly well.

References

- [1] Shumway, R.H. (1988) Applied Statistical Time Series Analysis. Prentice-Hall, Englewood Cliffs.
- [2] Hastie, T.J. and Tibshirani, T. (1993) Varying-Coefficient Models. *Journal of the Royal Statistical Society*, **55**, 757-796.
- [3] Hoover, D.R., Rice, J.A., Wu, C.O. and Yang, L.P. (1998) Nonparametric Smoothing Estimates of Time-Varying Coefficient Models with Longitudinal Data. *Biometrika*, **85**, 809-822. <http://dx.doi.org/10.1093/biomet/85.4.809>
- [4] Fan, J.Q. and Zhang, W.Y. (1999) Statistical Estimation in Varying-Coefficient Models. *Annals of Statistics*, **27**, 1491-1518. <http://dx.doi.org/10.1214/aos/1017939139>
- [5] Huang, J.Z., Wu, C.O. and Zhou, L. (2004) Polynomial Spline Estimation and Inference for Varying Coefficient Models with Longitudinal Data. *Statistica Sinica*, **14**, 763-788.

- [6] Thomas, D.C. (1981) General Relative-Risk Models for Survival Time and Matched Case-Control Analysis. *Biometrics*, **37**, 673-686. <http://dx.doi.org/10.2307/2530149>
- [7] Lustbader, E.D., Moolgavkar, S.H. and Venzon, D.J. (1984) Tests of the Null Hypothesis in Case-Control Studies. *Biometrics*, **40**, 1017-1024. <http://dx.doi.org/10.2307/2531152>
- [8] Kay, R. and Little, S. (1987) Transformations of the Explanatory Variables in the Logistic Regression Model for Binary Data. *Biometrika*, **74**, 495-501. <http://dx.doi.org/10.1093/biomet/74.3.495>
- [9] Qin, J. and Zhang, B. (1997) A Goodness-of-Fit Test for Logistic Regression Models Based on Case-Control Data. *Biometrika*, **84**, 609-618. <http://dx.doi.org/10.1093/biomet/84.3.609>
- [10] Qin, J., Berwick, M., Ashbolt, R., *et al.* (2002) Quantifying the Change of Melanoma Incidence by Breslow Thickness. *Biometrics*, **58**, 665-670. <http://dx.doi.org/10.1111/j.0006-341X.2002.00665.x>
- [11] Zhang, B. (2001) A Information Matrix Test for Logistic Regression Models Based on Case-Control Data. *Biometrika*, **88**, 921-932. <http://dx.doi.org/10.1093/biomet/88.4.921>
- [12] Zou, F., Fine, J.P. and Yandell, B.S. (2002) On Empirical Likelihood for a Semiparametric Mixture Model. *Biometrika*, **89**, 61-75. <http://dx.doi.org/10.1093/biomet/89.1.61>
- [13] White, I.R. and Thompson, S.G. (2003) Choice of Test for Comparing Two Groups, with Particular Application to Skewed Outcomes. *Statistics in Medicine*, **22**, 1205-1215.
- [14] Liu, X., Jiang, H. and Zhou, Y. (2013) Local Empirical Likelihood Inference for Varying-Coefficient Density-Ratio Models Based on Case-Control Data. *Journal of the American Statistical Association*, **109**, 635-646. <http://dx.doi.org/10.1080/01621459.2013.858629>
- [15] Thomas, D.R. and Grunkemeier, G.L. (1975) Confidence Interval Estimation of Survival Interval Estimation of Survival Probabilities for Censored Data. *Journal of the American Statistical Association*, **70**, 865-871. <http://dx.doi.org/10.1080/01621459.1975.10480315>
- [16] Owen, A. (2001) Empirical Likelihood. Chapman and Hall, New York. <http://dx.doi.org/10.1201/9781420036152>
- [17] Zhu, H.T., Ibrahim, J.G., Tang, N.S and Zhang, H. (2008) Diagnostic Measures for Empirical Likelihood of Generalized Estimating Equations. *Biometrika*, **95**, 489-507. <http://dx.doi.org/10.1093/biomet/asm094>
- [18] Xue, L. and Zhu, L. (2010) Empirical Likelihood in Nonparametric and Semiparametric Models. Science Press, Beijing.
- [19] Cook, R.D. and Weisberg, S. (1982) Residuals and Influence in Regression. Chapman and Hall, New York.
- [20] Wei, B., Lu, G. and Shi, J. (1990) Statistical Diagnostics. Publishing House of Southeast University, Nanjing.
- [21] Cai, Z., Fan, J. and Li, R. (2000) Efficient Estimation and Inferences for Varying-Coefficient Models. *Journal of American Statistical Association*, **95**, 888-902. <http://dx.doi.org/10.1080/01621459.2000.10474280>
- [22] Fan, J. and Zhang, W. (2008) Statistical Methods with Varying Coefficient Models. *Statistics and Its Interface*, **1**, 179-195. <http://dx.doi.org/10.4310/SII.2008.v1.n1.a15>
- [23] Prentice, R.L. and Pyke, R. (1979) Logistic Disease Incidence Models and Case-Control Studies. *Biometrika*, **66**, 403-411. <http://dx.doi.org/10.1093/biomet/66.3.403>