

## Surrogate Residuals for Discrete Choice Models

Journal:	<i>Journal of Computational and Graphical Statistics</i>
Manuscript ID	JCGS-19-067.R1
Manuscript Type:	Original Article
Keywords:	Categorical outcome, Model diagnostics, Multinomial logistic regression, Residual analysis

SCHOLARONE™  
Manuscripts

# Surrogate Residuals for Discrete Choice Models

Chao Cheng

Department of Mathematical Sciences, Tsinghua University

Rui Wang

Department of Public Economics, School of Economics, Xiamen University  
and

Heping Zhang\*

Department of Biostatistics, School of Public Health, Yale University

August 27, 2019

## Abstract

Discrete Choice Models (DCMs) are a class of models for modelling response variables that take values from a set of alternatives. Examples include logistic regression, probit regression, and multinomial logistic regression. These models are also referred together as generalized linear models. Although there exist methods for the goodness of fit of DCMs, defining intuitive residuals for such models has been difficult due to the fact that the responses are categorical values instead of continuous numbers. In this article, we propose the surrogate residual for DCMs based on the surrogate approach [Liu and Zhang, *J. Am. Stat. Assoc.*, **113**, 522 (2018) 845–854], which deals with an ordinal response. We consider categorical responses that may or may not be ordered. We shall show that our residual can be used to diagnose misspecification in the aspects of mean structure, individual-specific coefficients, and interaction effects.

*Keywords:* Categorical outcome; model diagnostics; multinomial logistic regression; residual analysis.

\*Address for correspondence: Heping Zhang, Department of Biostatistics, School of Public Health, Yale University, New Haven, CT, USA. Email: heping.zhang@yale.edu

# 1    Introduction

2  
3  
4  
5  
6 Discrete Choice Models (DCMs) (McFadden, 1978), a class of models designed for modelling  
7 the choice from a finite set of alternatives, are widely applied in medical research (Chan,  
8 2005), transportation demand forecasting (Ben-Akiva and Lerman, 1985), environmental  
9 analysis (Train, 1998), and evacuation modelling (Lovreglio et al., 2014). Specific forms  
10 of DCMs consist of logistic regression, probit regression, multinomial logistic regression,  
11 nested logit regression (Koppelman and Bhat, 2006) and so on. Plenty of diagnostic tests  
12 (Goeman and le Cessie, 2006; Pagan and Vella, 1989; Hausman and McFadden, 1981; Nagel  
13 and Hatzinger, 1992) were developed to check the validity of model assumptions for DCMs,  
14 whereas there are few useful tools for residual diagnosis for such models. Those diagnostic  
15 tests can provide useful measures for the probability to reject the null hypothesis, while  
16 the residual diagnosis can be more intuitive and revealing if appropriately developed. For  
17 instance, through the residual-by-covariate plot in the ordinary linear regression, we can  
18 straightforwardly review which parts of observed data are poorly fitted, thus producing  
19 potentially useful hints on the model improvement.

20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30 The difficulty of defining residuals for DCMs arises from the fact that the response  
31 variable is just a label of alternatives instead of a continuous value. Nagel and Hatzinger  
32 (1992) proposed a residual based on a quadratic-form statistic and applied the residual-by-  
33 leverage plot to help diagnose the model. However, this method results in an always-positive  
34 residual which may be inconvenient to interpret.

35  
36  
37  
38  
39 An intuitive approach is to assign numeric values for the discrete alternatives, and then  
40 applying the Pearson's residual, deviance residual or some other generalized residuals for  
41 above assigned numeric values (Dunn and Smyth, 1996; Gupta et al., 2008). However,  
42 this approach may be also problematic due to the following two reasons. First, how to  
43 assign the numeric values is arbitrary and subjective, lacking a sound rationale. Second,  
44 noting that the choice set is finite, the assigned numbers must be discrete; and the residuals  
45 based on discrete space could lead to unusual residual patterns even if the model is fitted  
46 correctly (Dunn and Smyth, 1996; Liu and Zhang, 2018). The following example illustrates  
47 the limitations of this approach.

1  
2  
3     **Example 1.** (Correct specification of the model). Suppose 2000 observations are generated  
4 from the following logistic regression  
5  
6

$$Pr(Y = 1|X) = \text{logit}(1 + 2X),$$

7  
8 where  $Y = \{1, 2\}$  is the response representing alternatives 1 and 2,  $X \sim \text{Unif}(-5, 5)$ , and  
9  
10  $\text{logit}(x) = \frac{e^x}{1+e^x}$ . Noting  $Y$  is a categorical variable, we may use a numeric value  $\tilde{Y} = \{1, 0\}$   
11 to represent the alternatives 1 and 2 respectively. Now, we use the true model to fit the  
12 simulated data, and obtain the Pearson's and deviance residuals, which are defined by  
13  
14  $r^{(P)} = (\tilde{y} - \hat{p})/\sqrt{\hat{p}(1 - \hat{p})}$  and  $r^{(D)} = \mathbb{I}(\tilde{y} - \hat{p})\sqrt{2[(1 - \tilde{y})\ln(1 - \hat{p}) - \tilde{y}\ln(\hat{p})]}$ , respectively.  
15  
16 Here,  $\hat{p} = \hat{P}(\tilde{Y} = 1|X = x)$ , and  $\mathbb{I}(c) = 1$  if  $c > 0$  and 0 otherwise.  
17  
18

19  
20     Figure 1 (panels a - d) presents the residual-by-covariate plots ( $r$ -by- $x$ ) and quantile-  
21 quantile (QQ) plots (the empirical distribution of  $r$  versus the standard normal distribution)  
22 for both Pearson's and deviance residuals. It is difficult to observe anything useful from  
23 Pearson's and deviance residuals even if the logistic model is correctly fitted. In fact, how  
24 these residuals are defined is problematic. First, the distribution of the residual variable  
25 conditional on covariate, i.e  $R|X$ , varies across the values of  $X$  (see Figure 1 (a) and (c)).  
26 Second, the null distribution of  $R$  may depend on the distribution of  $X$  (see Figure 1 (b)  
27 and (d)). These issues limit the use of these residuals for model diagnostic.  
28  
29

30  
31     Liu and Zhang (2018) provided the surrogate residual for ordinal regression. Specifi-  
32 cally, they constructed a continuous variable,  $S$ , named surrogate variable, which is given  
33 conditionally on the ordinal outcome variable  $Y$ , and then defined the surrogate residual  
34 as the difference of  $S$  and its expectation under the null hypothesis. Due to the continuity  
35 of  $S$ , the surrogate residual avoids the problems encountered in example 1. This advantage  
36 motivated us to develop a residual based on a continuous surrogate variable rather than  
37 the discrete outcomes.  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

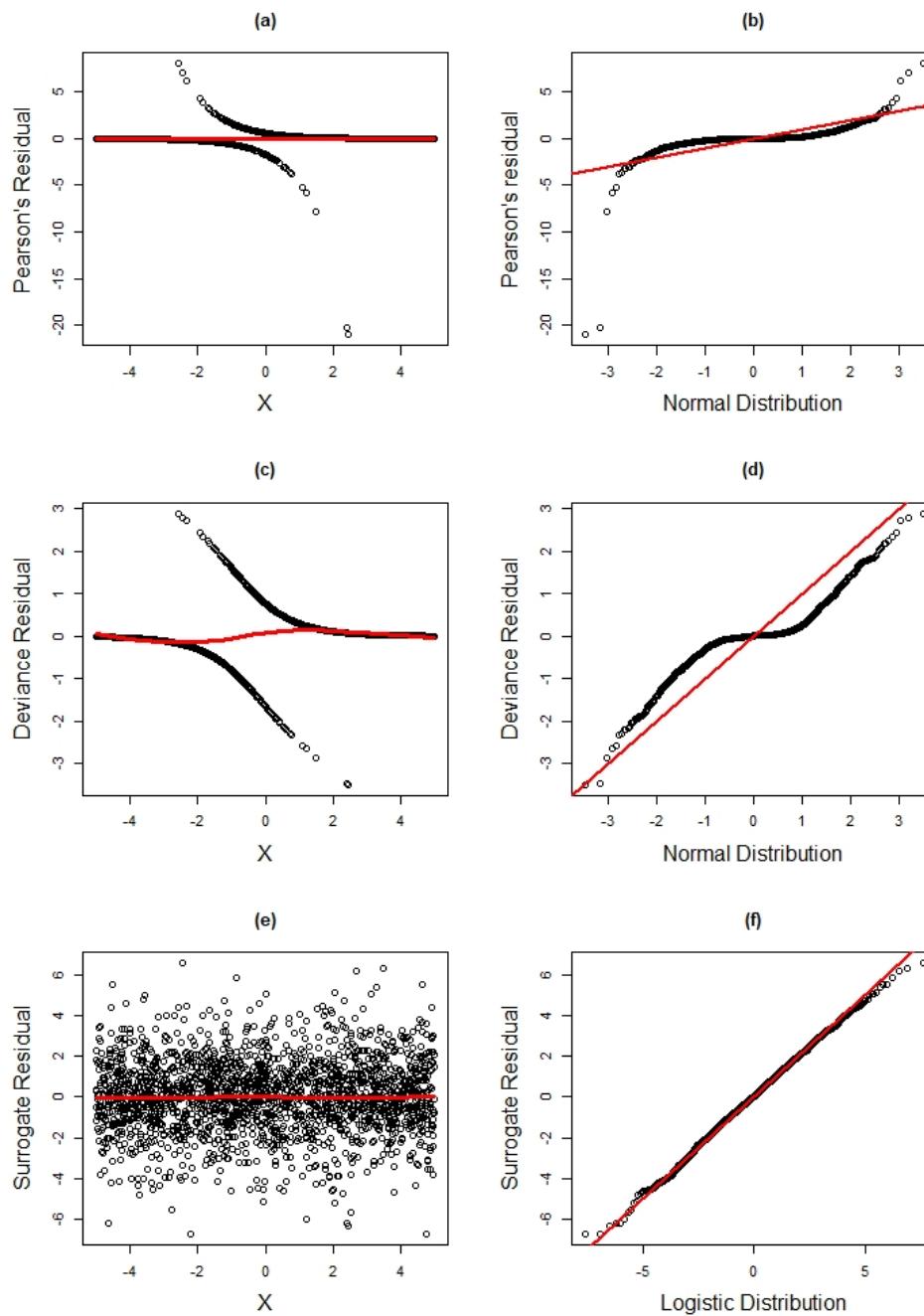


Figure 1: Model diagnostics using the Pearson's residual (upper row), deviance residual (middle row) and surrogate residual (lower row) in the case that the logistic regression is correctly fitted. Figures (a), (c) and (e) present the plots of the residual-by-covariate ( $r$ -by- $x$ ), and loess curves (red solid) are added as references. Figures (b), (d) and (f) are the corresponding QQ-plots (quantile-quantile plots) of above mentioned residuals.

In this article, we are concerned with DCMs for which the response variable takes discrete choices without orders. This response variable is different from the ordinal response investigated by Liu and Zhang (2018). However, we extend the use of the surrogate approach to define the residual for DCMs. Because discrete choices have no order, each alternative has its unique utility function. This is the major difference between our work and that of Liu and Zhang (2018). The multiple utility functions create additional challenges in defining the residual. First, an one-dimensional residual is no longer appropriate because it cannot represent the multi-dimensional utilities. Second, the residual needs to reflect the relationship among different alternatives. In this work, the surrogate residual  $\mathbf{R}$  for DCMs is a vector for which each dimension corresponding to one particular alternative, and has a reference distribution with mean  $\mathbf{0}$  and is independent of the covariate  $\mathbf{X}$  if the assumed model holds. Based on theoretical and graphical analyses, we shall show that, by comparing the surrogate residual with its reference distribution and synthesizing the information across different residual dimensions, our residual is able to detect various misspecifications, including detecting the misspecified mean structure, checking individual-specific coefficients, and detecting interaction effects.

The framework of this paper is as follows: Discrete choice models are briefly introduced in Section 2. In Section 3, we define the surrogate residual for DCMs, and present its theoretical and graphical properties. We generalize the surrogate residual for general statistical models with unordered categorical responses in Section 4. In Section 5, we re-analyze a publicly available dataset (Hoogendoorn-Lanser et al., 2015) to demonstrate the use of the surrogate residual in real data analysis. We conclude this paper with a brief discussion in Section 6.

## 2 Discrete Choice Models

Suppose a categorical variable  $Y$  that has  $J$  alternatives  $\{1, \dots, J\}$ . Consider there are  $J$  utility levels

$$U_j = V_j(\mathbf{X}_j, \boldsymbol{\beta}_j) + \epsilon_j \quad (1)$$

corresponding to alternative  $j = 1, \dots, J$ , where  $V_j(\mathbf{X}_j, \boldsymbol{\beta}_j)$  is the deterministic part depending on covariates  $\mathbf{X}_j = [X_{j1}, \dots, X_{jp}]^T$  and unknown coefficients  $\boldsymbol{\beta}_j$ , and  $\epsilon_j$  is the random part. The random parts  $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_J]^T$  are assumed to follow a  $J$ -dimensional continuous distribution  $G(\boldsymbol{\epsilon})$ . In a DCM,  $Y = y$  if and only if  $\forall j \neq y, U_y > U_j$ ; i.e, the probability of selecting alternative  $y$  is

$$\Pr(Y = y) = \Pr(U_y > U_j, \forall j \neq y), \quad (2)$$

where  $y \in \{1, 2, \dots, J\}$ . In order to ensure the uniqueness of the unknown parameters  $\boldsymbol{\beta}$ , some specific models select an alternative as reference and force its whole utility (both deterministic and random parts) as zero. Without loss of generality, we select the last alternative  $J$  as the reference.

The following are a few examples for how model (1) can be specified; for DCMs whose choice set contains 2 alternatives, we have:

- Logistic regression, where  $U_2 \equiv 0$  and  $\epsilon_1$  follows a logistic distribution with c.d.f  $F(\gamma) = e^\gamma / (1 + e^\gamma)$ .
- Probit regression, where  $U_2 \equiv 0$  and  $\epsilon_1$  follows a standard normal distribution.

Below are multinomial choice models whose choice set contains 3 or more alternatives:

- Multinomial logistic regression, where  $\epsilon_j$ 's are mutually independent, and follow a standard gumbel distribution with c.d.f  $F(\gamma) = \exp(-e^{-\gamma})$ .
- Multinomial probit regression, where  $U_J \equiv 0$  and  $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_{J-1}]^T \sim N(0, \Sigma_{(J-1) \times (J-1)})$ .
- Nested logit regression, where the c.d.f of  $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_J]^T$  is  $F(\epsilon_1, \dots, \epsilon_J) = \exp \left( - \sum_{s=1}^S \left( \sum_{j \in B_s} e^{-\epsilon_j / \lambda_s} \right)^{\lambda_s} \right)$ .

**Remark 1.** Gumbel distribution will be used later, so we briefly describe its known properties here. The c.d.f of Gumbel distribution is  $F(x; \theta, \alpha) = \exp(-e^{(x-\theta)/\alpha})$ , where  $\theta$  is the location parameter and  $\alpha$  is the scale parameter. The variance is  $\frac{\pi^2}{6}\alpha^2$ , and the mean is given by  $\theta + \zeta\alpha$ , where  $\zeta \approx 0.5772$  denotes the Euler-Mascheroni constant. Moreover, it is a positive skewed distribution with the moment coefficient of skewness 1.1423. The standard

Gumbel distribution is the Gumbel distribution with  $\theta = 0$  and  $\alpha = 1$ , whose mean is  $\zeta$  and variance is approximately equal to 1.6449.

### 3 Surrogate Residual for DCMs

#### 3.1 Definition

In this section, we propose a residual for DCMs based on the surrogate approach proposed by Liu and Zhang (2018). We shall (1) find a hypothetical variable  $\mathbf{Z}$  following a continuous distribution, (2) construct a joint distribution  $(\mathbf{Z}, Y)$  to link  $\mathbf{Z}$  and  $Y$  together, and (3) define a surrogate variable  $\mathbf{S}$  as the conditional distribution of  $\mathbf{Z}$  given  $Y$ . Finally, the surrogate residual  $\mathbf{R}$  is the difference between  $\mathbf{S}$  and the null distribution of  $\mathbf{S}$  conditional on  $\mathbf{X}$ , i.e  $E_0(\mathbf{S}|\mathbf{X})$ .

In DCMs, the utility vector  $\mathbf{U} = [U_1, \dots, U_J]^T$  is a convenient hypothetical variable because  $\mathbf{U}$  is a continuous variable due to the continuity of  $\epsilon$ ; and importantly, equation (2) induces a natural joint distribution linking  $\mathbf{U}$  and  $Y$ . Specifically, the joint distribution can be set as

$$Y \triangleq y \iff U_y > U_j \quad \forall j \neq y.$$

Now, following the procedure of the surrogate approach, we let the surrogate variable  $\mathbf{S}$  be a  $J$ -dimensional random variable following the conditional distribution of  $\mathbf{U}$  given the observation  $Y$ . More specifically,  $\mathbf{S} \sim \mathbf{U}|(Y, \mathbf{X})$  follows a truncated distribution obtained by truncating the distribution of  $\mathbf{U}|\mathbf{X} = \mathbf{V}(\mathbf{X}, \boldsymbol{\beta}) + \boldsymbol{\epsilon}$  using the region  $U_Y > U_j$  ( $\forall j \neq Y$ ), where  $\mathbf{U} = [U_1, \dots, U_J]^T$  and  $\mathbf{V}(\mathbf{X}, \boldsymbol{\beta}) = [V_1(\mathbf{X}_1, \boldsymbol{\beta}_1), \dots, V_J(\mathbf{X}_J, \boldsymbol{\beta}_J)]^T$ . We name  $S_j$  ( $j \in \{1, \dots, J\}$ ) as the surrogate response for alternative  $j$  or utility  $U_j$ . We define the surrogate residual  $\mathbf{R}$  as

$$\mathbf{R} = \mathbf{S} - E_0(\mathbf{S}|\mathbf{X}) = \mathbf{S} - E(\mathbf{U}|\mathbf{X}) = \mathbf{S} - \mathbf{V}(\mathbf{X}, \boldsymbol{\beta}) - \int \boldsymbol{\epsilon} dG(\boldsymbol{\epsilon}),$$

where  $E_0(\mathbf{S}|\mathbf{X})$  denotes the conditional mean of  $\mathbf{S}$  given  $\mathbf{X}$  under the null hypothesis that the model is specified correctly. In practice, given the observed data  $\{y, \mathbf{x}\}$  and a fitted model, we estimate the conditional distribution  $\mathbf{S} \sim \mathbf{U}|(y, \mathbf{x})$  by plugging in the estimated

parameter  $\hat{\beta}$ , and randomly draw a sample  $s$  from this plug-in distribution. Finally, the residual  $r$  is given by  $r = s - V(x, \hat{\beta}) - \int \epsilon dG(\epsilon)$ .

We should note that  $r$  is a sample from  $\hat{R}_{\hat{\beta}}$  rather than a realization of  $R \equiv R_{\beta}$ , but we can show that when  $\hat{\beta} \rightarrow \beta$  in probability,  $\hat{R}_{\hat{\beta}}$  converges to  $R_{\beta}$  in distribution. In addition, the properties of  $R_{\beta}$  apply to  $\hat{R}_{\hat{\beta}}$  asymptotically (see Theorems 3 and 4).

In practice, we need to draw a sample from  $S$ , which follows a truncated  $J$ -dimensional distribution. As we know, *Inverse CDF Sampling* and *Rejection Sampling* are two effective methods to generate samples in one dimension. Due to the multidimensionality of  $S$ , directly applying either method is quite difficult. In this paper, we introduced the Gibbs sampling method to generate this truncated multi-dimensional distribution by transforming the multidimensional sampling problem into a series of full conditional distribution (which is one-dimensional) sampling issues. Notice that  $S$  is a distribution subject to a specific choice, and hence the full conditional distribution of each dimension can be simplified and generated by *Inverse CDF Sampling* or any other regular univariate sampling algorithm. The specific procedure of generating samples from  $S$  is deferred to Appendix A in the supplementary material.

**Remark 2.** As discussed in Section 2, some DCMs force  $U_J \equiv 0$ , and assume  $\epsilon$  follows a  $(J-1)$ -dimensional distribution  $G(\epsilon_1, \dots, \epsilon_{J-1})$ . In this case, the c.d.f of the surrogate variable  $S$  follows a truncated distribution obtained by truncating the distribution of  $U|X = V(X, \beta) + \epsilon$  using the region  $U_j > U_j$ , where  $V(X, \beta) = [V_1(X_1, \beta_1), \dots, V_{J-1}(X_{J-1}, \beta_{J-1})]^T$ ,  $U = [U_1, \dots, U_{J-1}]^T$  and  $U_J$  is fixed as 0. This distribution is a special form of distribution of  $S$  above, in which  $U_J \equiv 0$  can be seen as a degenerate distribution. Thus, the properties derived in Section 3.2 are still applicable in this case.

## 3.2 Theoretical Properties

We investigate the theoretical properties of the surrogate variable  $S$  and the residual variable  $R$  in this subsection. These properties provide a theoretical foundation for model diagnostics. Firstly, suppose the true utility levels are

$$\tilde{U}_j = \tilde{V}_j(X_j, \tilde{\beta}_j) + \tilde{\epsilon}_j \quad j = 1, \dots, J,$$

1  
2  
3 where  $\tilde{\boldsymbol{\epsilon}} = [\tilde{\epsilon}_1, \dots, \tilde{\epsilon}_J]^T \sim \tilde{G}(\tilde{\boldsymbol{\epsilon}})$ , and the true model for  $Y$  is  
4  
5

$$Pr(Y = y) = Pr(\tilde{U}_y > \tilde{U}_j, \forall j \neq y), \quad (3)$$

6  
7 where  $y = \{1, \dots, J\}$ . Then, the distribution of  $\mathbf{S}$  is  
8  
9

$$\begin{aligned} Pr(\mathbf{S} \leq \mathbf{c}) &= \sum_{y=1}^J Pr(\mathbf{S} \leq \mathbf{c} | Y = y) Pr(Y = y) \\ &= \sum_{y=1}^J \frac{Pr(\mathbf{U} \leq \mathbf{c} \wedge U_y > U_j, \forall j \neq y)}{Pr(U_y > U_j, \forall j \neq y)} Pr(\tilde{U}_y > \tilde{U}_j, \forall j \neq y), \end{aligned} \quad (4)$$

10  
11 where  $\mathbf{U} = [U_1, \dots, U_J]^T$  is the assumed utility levels in equation (1),  $\tilde{\mathbf{U}} = [\tilde{U}_1, \dots, \tilde{U}_J]^T$   
12 is the true utility levels, and  $\mathbf{c} = [c_1, \dots, c_J]^T$  is an arbitrary but fixed point in  $\mathbb{R}^J$ . Then  
13 denote the area  $(-\infty, \mathbf{c})$  by  $\mathcal{C}$ ,  $\{U_y > U_j, \forall j \neq y\}$  by  $\mathcal{D}_y$ , and  $\{\tilde{U}_y > \tilde{U}_j, \forall j \neq y\}$  by  $\tilde{\mathcal{D}}_y$ .

14 Now, the distribution (4) can be rewritten as  
15  
16

$$Pr(\mathbf{S} \leq \mathbf{c}) = \sum_{y=1}^J \frac{Pr(\mathbf{U} \in \mathcal{C} \cap \mathcal{D}_y)}{Pr(\mathbf{U} \in \mathcal{D}_y)} Pr(\tilde{\mathbf{U}} \in \tilde{\mathcal{D}}_y). \quad (5)$$

17  
18 Equation (5) shows that the distribution of  $\mathbf{S}$  is determined jointly by the assumed and the  
19 true model for  $Y$ . If the assumed model agrees with the true model for  $Y$ , i.e.  $\boldsymbol{\beta} = \tilde{\boldsymbol{\beta}}$ ,  $\mathbf{V} =$   
20  $\tilde{\mathbf{V}}$ ,  $G = \tilde{G}$ , we have  $\mathbf{U} = \tilde{\mathbf{U}}$  and  $\mathcal{D}_y = \tilde{\mathcal{D}}_y$  for all  $y = 1, \dots, J$ . Furthermore, noticing that  
21  $\{\mathcal{D}_1, \dots, \mathcal{D}_J\}$  is a partition of  $\mathbb{R}^J$ , we have  
22

$$\begin{aligned} Pr(\mathbf{S} \leq \mathbf{c}) &= \sum_{y=1}^J \frac{Pr(\mathbf{U} \in \mathcal{C} \cap \mathcal{D}_y)}{Pr(\tilde{\mathbf{U}} \in \tilde{\mathcal{D}}_y)} Pr(\tilde{\mathbf{U}} \in \tilde{\mathcal{D}}_y) \\ &= \sum_{y=1}^J Pr(\mathbf{U} \in \mathcal{C} \cap \mathcal{D}_y) = Pr\left(\mathbf{U} \in \bigcup_{y=1}^J (\mathcal{C} \cap \mathcal{D}_y)\right) \\ &= Pr\left(\mathbf{U} \in \mathcal{C} \cap (\bigcup_{y=1}^J \mathcal{D}_y)\right) = Pr(\mathbf{U} \leq \mathbf{c}), \end{aligned} \quad (6)$$

23  
24 i.e., the surrogate variable  $\mathbf{S}$  has the same distribution as  $\mathbf{U}$ . This equation immediately  
25 yields the following theorem.  
26  
27

28  
29 **Theorem 1.** *If the assumed model (2) agrees with the true model (3) (i.e.,  $\boldsymbol{\beta} = \tilde{\boldsymbol{\beta}}$ ,  $\mathbf{V} =$   
30  $\tilde{\mathbf{V}}$ ,  $G = \tilde{G}$ ), then the following results hold*

- 31  
32 (1) *The surrogate variable  $\mathbf{S}$  has the same distribution as  $\mathbf{U}$ , i.e.  $\mathbf{S}|\mathbf{X} = \mathbf{U}|\mathbf{X} =$   
33  $\mathbf{V}(\mathbf{X}, \boldsymbol{\beta}) + \boldsymbol{\epsilon}$ .*
- 34  
35 (2) *The residual variable  $\mathbf{R}$ , independent from  $X$ , follows the distribution  $G(\mathbf{r} + \int \boldsymbol{\epsilon} dG(\boldsymbol{\epsilon}))$ ,*  
36 *i.e.,  $Pr(\mathbf{R} \leq \mathbf{r}) = Pr(\mathbf{R} \leq \mathbf{r}|\mathbf{X}) = G(\mathbf{r} + \int \boldsymbol{\epsilon} dG(\boldsymbol{\epsilon}))$ .*

37  
38 Theorem 1 leads to the following result for checking the assumed model.  
39  
40

**Theorem 2.** If the assumed model (2) agrees with the true model (3) (i.e.,  $\beta = \tilde{\beta}, V = \tilde{V}, G = \tilde{G}$ ), then the residual variable  $R$  has the following properties

- (1) (Zero mean)  $E(R|X) = 0$ .
- (2) (Homogeneous variance)  $Var(R|X) = Var(\epsilon)$  is a constant matrix, not depending on  $X$ .

Theorems 1 and 2 establish the theoretical foundation of using  $R$  for model diagnosis. Our surrogate residual  $R$  has similar properties to the residual in the ordinary linear regression. By plotting the quantile-quantile plot of  $r$  with the reference distribution  $G(r + \int \epsilon dG(\epsilon))$ , we can assess whether  $R$  agrees with  $\tilde{R}$ . By drawing the residual-by-covariate plot, we can evaluate what aspects of our assumed models are poorly fitted.

Theorems 1 and 2 concern the surrogate residual variable  $R \equiv R_{\beta}$ . Theorems 3 and 4 present the parallel results when  $\hat{\beta} = \beta + o_p(1)$ , here  $o_p(1)$  denotes a term vanishing to zero in probability.

**Theorem 3.** If the assumed model (2) agrees with the true model (3) (i.e.,  $\beta = \tilde{\beta}, V = \tilde{V}, G = \tilde{G}$ ), then  $\hat{R}_{\hat{\beta}}$  converges to  $G(r + \int \epsilon dG(\epsilon))$  in distribution as  $n \rightarrow \infty$ , conditional or unconditional on  $X$ , i.e.,  $Pr(\hat{R}_{\hat{\beta}} \leq r) = Pr(\hat{R}_{\hat{\beta}} \leq r|X) = G(r + \int \epsilon dG(\epsilon)) + o(1)$ .

**Theorem 4.** If the assumed model (2) agrees with the true model (3) (i.e.,  $\beta = \tilde{\beta}, V = \tilde{V}, G = \tilde{G}$ ), then  $\hat{R}_{\hat{\beta}}$  has the following properties

- (1)  $E(\hat{R}_{\hat{\beta}}|X) = o(1)$ .
- (2)  $Var(\hat{R}_{\hat{\beta}}|X)$  is a constant matrix, not depending on  $X$ , except a vanishing term  $o(1)$ .

Appendix B in the supplementary material provides detailed proofs for Theorems 3 and 4.

### 3.3 Graphical Properties

In this subsection, using several numerical examples, we examine graphical properties of our proposed residual, when the model is specified correctly or misspecified. We first revisit Example 1.

1  
2  
3     **Example 1.** (Continued) We choose alternative 2 as reference; thus, the true utilities are  
4      $U_1 = 1 + 2X + \epsilon_1$  and  $U_2 \equiv 0$ , corresponding to alternative 1 and 2 respectively. Here,  $\epsilon_1$   
5     follows a logistic distribution. When the model is specified correctly, our surrogate residual  
6     is an one-dimensional random variable following a logistic distribution. We obtain the sur-  
7     rogate residual, and draw the residual-by-covariate plot and QQ plot in Figure 1 (e) and  
8     (f). As the residual-by-covariate plot shows, our surrogate residuals are randomly fluctu-  
9     ating around zero and independent from the covariate  $X$ , exhibiting no unusual pattern.  
10    Based on the QQ plot, we can see that the residuals are in accordance with the reference  
11    distribution.  
12  
13

14  
15    We should note that Example 1 is a very special instance of DCMs as it has two  
16    alternatives only. Most DCMs, however, contain  $J \geq 3$  alternatives, indicating that the  
17    surrogate residual is a  $J$  or  $J - 1$  (if  $U_J \equiv 0$  is constrained) dimensional vector. Unlike  
18    the residuals in ordinary regression and other generalized linear models, this characteristic  
19    creates additional difficulties in model diagnostics. First, we need to evaluate whether  
20    every dimension of the residual vector satisfies the null hypothesis or not. Second, we need  
21    to synthesize the information across different residual dimensions and identify which part  
22    of the model assumptions is misspecified.  
23  
24

25    In the following, we use the multinomial logistic regression with a 3-alternative response  
26    as examples to illustrate how to detect a misspecified mean structure, check the propor-  
27    tionality assumption, and inspect interaction effects. Similar steps are applicable to general  
28    DCMs.  
29  
30

### 31    3.3.1 Detecting Misspecified Mean Structure 32

33    After fitting a discrete choice model it is important to determine whether the mean struc-  
34    tures (i.e the deterministic parts  $V_j(\mathbf{X}_j, \boldsymbol{\beta}_j)$ ) of all utilities are correctly specified before  
35    performing inference.  
36  
37

38    We can detect the misspecified mean structure based on the residual-by-covariate plots.  
39    If the assumed model holds, our surrogate residual should be randomly distributed around  
40    zero, independent from  $X$ ; and if the mean structure is misspecified, the patterns of our  
41    residuals will deviate from the expected behavior.  
42  
43

surrogate residual depend on  $X$ .

**Example 2.** We generate 2000 observations from a multinomial logistic regression with the following utility levels

$$U_j = \beta_{j0} + \beta_{j1}X + \beta_{j2}X^2 + \epsilon_j, \quad j = 1, 2, 3,$$

where  $\epsilon_j$  follows a standard Gumbel distribution;  $X$  is an individual-specific covariate generated from  $Unif(-5, 5)$ ;  $\beta_{30} = \beta_{31} = \beta_{32} = 0$  as the reference;  $\beta_{10} = 1$ ,  $\beta_{11} = 7$ ,  $\beta_{12} = -2$ ; and  $\beta_{20} = 8$ ,  $\beta_{21} = 2$ ,  $\beta_{22} = -4$ . In order to examine diagnostic power of our residuals when the mean structure is misspecified, we do not include the quadratic term  $X^2$  in the assumed model. Instead, we fit the model with only a linear term of  $X$ , i.e we assume the utilities  $U_j = \beta_{j0} + \beta_{j1}X + \epsilon_j$ , for  $j = 1, 2, 3$ . We obtain the surrogate residual  $\mathbf{r} = [r_1, r_2, r_3]^T$  corresponding to the utilities 1, 2 and 3. The residual-by-covariate plots and QQ-plots of our surrogate residuals are shown in the left and right columns of Figure 2.

All the residual-by-covariate plots exhibit wave-like shapes, suggesting a misspecified mean structure. Specifically, we can see  $r_1$  slants upward for  $X \in (1, 3.5)$  and goes downward when  $X > 3.5$ . Similarly,  $r_2$  seems to be positive for  $X \in (-1, 1)$ ;  $r_3$  negative when  $X \in (-1, 3.5)$  and positive for  $X > 3.5$ . These patterns suggest that a linear form of  $X$  in the mean structure may not be enough, and a higher order term of  $X$  is needed. For comparison, we provide the residual-by-covariate plots of Pearson's and deviance residuals in Figure 3. Here, the Pearson's residual for multinomial logistic regression is given by  $\mathbf{r}^{(P)} = [r_1^{(P)}, \dots, r_J^{(P)}]$ , with

$$r_j^{(P)} = (\tilde{y}_j - \hat{p}_j) / \sqrt{\hat{p}_j(1 - \hat{p}_j)}, \quad j = 1, \dots, J,$$

and deviance residual  $\mathbf{r}^{(D)} = [r_1^{(D)}, \dots, r_J^{(D)}]$ , with

$$r_j^{(D)} = \mathbb{I}(\tilde{y}_j - \hat{p}_j) \sqrt{2[(1 - \tilde{y}_j) \ln(1 - \hat{p}_j) - \tilde{y}_j \ln(\hat{p}_j)]}, \quad j = 1, \dots, J,$$

where  $\hat{p}_j$  is the estimated possibility of  $Y = j$ , and  $\tilde{y}_j = 1$  if  $Y = j$  and 0 otherwise. Both definitions are natural extension of the Pearson's and deviance residuals for the binary

logistic regression. Similar to the surrogate residual, both are  $J$ -dimensional vectors corresponding to  $U_1$  to  $U_J$ . As shown in Figure 3, both  $\mathbf{r}^{(P)}$  and  $\mathbf{r}^{(D)}$  detect a wave-like shape but indicate three separate parts. As a next step, we added a quadratic term of  $X$ . The updated surrogate residual plots are displayed in Figure S1 in the supplementary material. We can see that all plots exhibit no abnormal pattern.

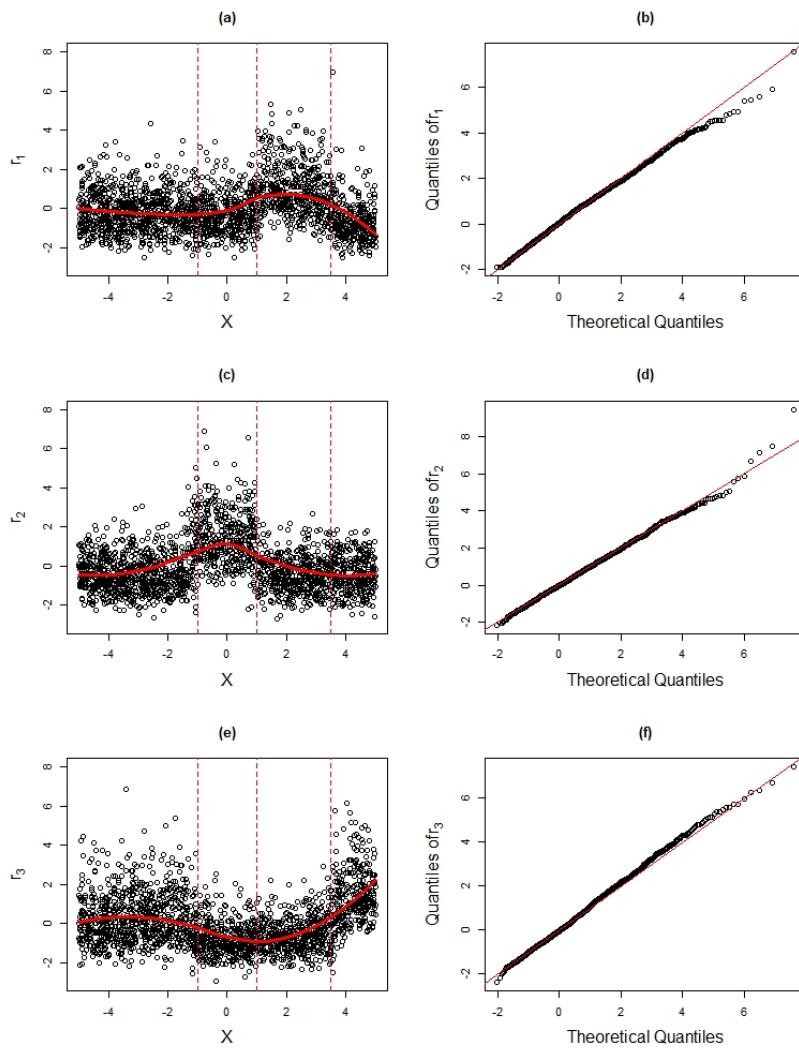


Figure 2: Model diagnostics when the mean structure in the assumed model is misspecified. Figures (a), (c) and (e) show the residuals-by-covariate plots of  $r_1$ ,  $r_2$  and  $r_3$  respectively; vertical dashed lines  $X = -1$ ,  $1$ , and  $3.5$ , and loess curves (red solid) are added as references. Figures (b), (d) and (f) are the corresponding QQ-plots (quantile-quantile plots) of above mentioned residuals.

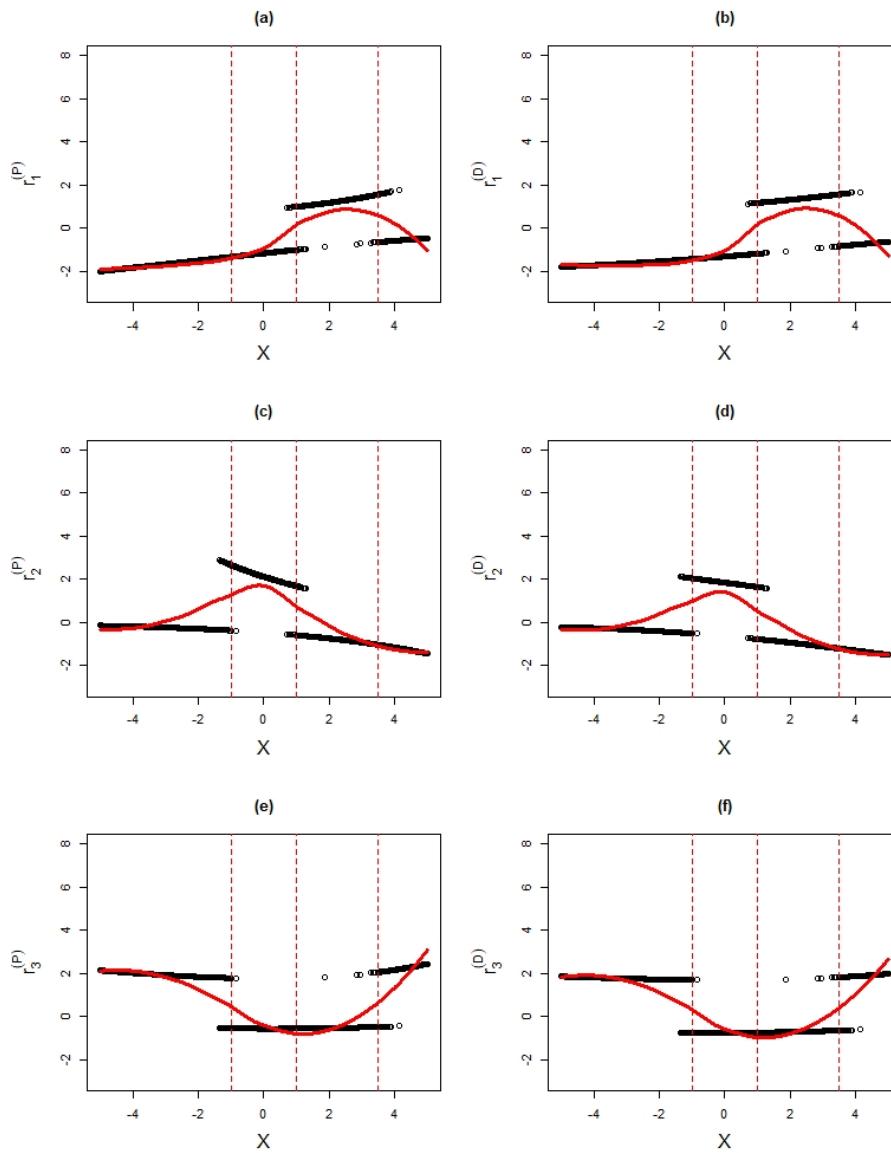


Figure 3: Model diagnostics when the mean structure in the assumed model is misspecified. The left and right columns display the residuals-by-covariate plots of Pearson's ( $r_j^{(P)}$ ) and deviance ( $r_j^{(D)}$ ) residuals, respectively. Vertical dashed lines  $X = -1, 1$ , and  $3.5$ , and loess curves (red solid) are added as references.

### 3.3.2 Checking the individual-specific coefficients

The covariates and coefficients in the DCMs can be “individual-specific”, i.e they are identical among all alternatives, or “alternative-specific”, i.e they are different among al-

ternatives  $1, \dots, J$ . Many empirical analyses of DCMs assume the coefficients of some alternative-specific variable are individual-specific, i.e.,  $\beta_{1k} = \beta_{2k} = \dots = \beta_{Jk} =: \beta_k$  for alternative-specific variable  $\tilde{\mathbf{X}}_k = [X_{1k}, \dots, X_{Jk}]$ . Such an assumption is adopted widely in practice to achieve parsimonious models. We demonstrate in the example below that our surrogate residual offers a simple way to check this assumption.

**Example 3.** We generated 2000 observations from a multinomial regression model with the following utility levels

$$U_j = \beta_{j0} + \beta_{j1}X_{j1} + \epsilon_j, \quad j = 1, 2, 3,$$

where  $\epsilon_j$  follows a standard Gumbel distribution;  $\mathbf{X}_1 = [X_{11}, X_{21}, X_{31}]^T$  are alternative-specific covariates independently generated from  $Unif(-5, 5)$ ;  $\beta_{10} = \beta_{20} = \beta_{30} = 0$  are the intercepts;  $\beta_{11} = 10$ ,  $\beta_{21} = 5$ ,  $\beta_{31} = -10$  are alternative-specific coefficients. It is of interest to check if it is reasonable to assume  $\beta_{11} = \beta_{21} = \beta_{31} =: \beta_1$ . Now, we fit the observations based on an individual-specific coefficient, i.e.,  $U_j = \beta_{j0} + \beta_1X_{j1} + \epsilon_j$  for  $j = 1, 2, 3$ .

If the true model contains an alternative-specific coefficient but the fitted model forces an individual-specific coefficient, we can detect this misspecification using the surrogate residuals. We simply draw  $r_j$ -by- $x_{j1}$  plots, for  $j = 1, 2, 3$ , and detect if there is any trend inside these three plots. Or more intuitively, we simply draw  $s_j$ -by- $x_{j1}$  plots, for  $j = 1, 2, 3$ , and compare the slopes. The slopes of  $s_j$  by  $x_{j1}$  for all  $j = 1, 2, 3$  should be the same, or imply an alternative-specific  $\beta_1$ . Figure S2 in the supplementary material displays the surrogate variable versus covariate plots, which indicates a positive slope for  $s_1$ -by- $x_{11}$  and  $s_2$ -by- $x_{21}$  but a negative trend for  $s_3$ -by- $x_{31}$ . Thus, an alternative-specific coefficient is detected. We also provide the surrogate-by-covariate plots of  $\mathbf{r}^{(P)}$  and  $\mathbf{r}^{(D)}$  in Figure S3 in the supplementary material. Here, the surrogate response for  $\mathbf{r}^{(P)}$  and  $\mathbf{r}^{(D)}$  are given by  $s_j^{(P)} = V_j(\mathbf{X}_j, \hat{\beta}_j) + r_j^{(P)}$  and  $s_j^{(D)} = V_j(\mathbf{X}_j, \hat{\beta}_j) + r_j^{(D)}$ ,  $j = 1, \dots, J$ , respectively. All the scatter points of  $s_j^{(P)}$  and  $s_j^{(D)}$  have a clear positive trend but it is hard to detect the misspecified alternative-specific coefficient as they cluster in two strips.

### 3.3.3 Detecting Interaction Effects

Another common issue in model development is to determine if there are important interactions between the covariates. We illustrate how to check whether or not such interaction effects are missing from the assumed model.

**Example 4.** We generated 2000 observations from a multinomial regression model with the following utility levels

$$U_j = \beta_{j0} + \beta_{j1}X_1 + \beta_{j2}X_2 + \beta_{j3}X_1X_2 + \epsilon_j, \quad j = 1, 2, 3,$$

where  $\epsilon_j$  follows a standard Gumbel distribution;  $X_1 \sim Unif(-5, 5)$  and  $X_2 \sim Bernoulli(0.5)$  are individual-specific variables;  $\beta_{30} = \beta_{31} = \beta_{32} = \beta_{33} = 0$  are references;  $\beta_{10} = -10$ ,  $\beta_{11} = -5$ ,  $\beta_{12} = 5$ ,  $\beta_{13} = 3$ ;  $\beta_{20} = -10$ ,  $\beta_{21} = 5$ ,  $\beta_{22} = -5$ ,  $\beta_{23} = 6$ . The indicator variable  $X_2$  can be seen as a factor with levels “Treatment” and “Control.”

To check if there are any interaction effects between  $X_1$  and  $X_2$ , we just fit two models: The first model corresponds to the simulated control group, while the second model corresponds to the treatment group. Both models assume  $U_j = \beta_{j0} + \beta_{j1}X_1 + \epsilon_j$  for  $j = 1, 2, 3$ . The surrogate response versus  $x_1$  plots can help us detect if there are any interaction effects between  $X_1$  and  $X_2$ . If the true model does not contain an interaction term  $X_1X_2$ , the slopes of  $s_j$  ( $j = 1, 2, 3$ ) by  $x_1$  between the control and treatment groups should be exactly the same, otherwise indicating an interaction effect term between  $X_1$  and  $X_2$ . In our surrogate-by-covariates plots (see Figure S4 in the supplementary material), the control group displays a stronger negative association between  $s_1$  and  $x_1$  comparing with the treatment group; but the treatment group shows a stronger positive association between  $s_2$  and  $x_1$  comparing with the control group. Therefore, the interaction effect should be considered based on our surrogate residuals. As shown in Figures S5 and S6 in the supplementary material, the surrogate responses of Pearson’s and deviance residuals exhibit strange patterns and cannot be used to detect interaction terms.

## 3.4 Practical Remarks

In this subsection, we give two practical remarks related to our surrogate residuals.

### 3.4.1 Normalized Surrogate Residual

From Section 3.2, we see that the null distribution of our surrogate residual depends on the distribution of the random part  $\epsilon$ . Specifically,  $\mathbf{R} \sim G(\mathbf{R} + \int \epsilon dG(\epsilon))$  if the assumed model holds. Lack of knowledge for the distribution  $G(\cdot)$  may lead to difficulties in practice. For instance, the null distribution of the surrogate residual for multinomial logistic regression is a Gumbel distribution, which is slightly positive skewed. Our inference may be biased if this positive-skewness is not considered. Here, for the convenience of the residual diagnosis, we may normalize the surrogate residuals to the standard normal distribution. Specifically, given the surrogate residual  $\mathbf{r} = [r_1, \dots, r_J]^T$ , the normalized surrogate residual  $\mathbf{r}^{(N)} = [r_1^{(N)}, \dots, r_J^{(N)}]^T$  can be obtained by

$$r_j^{(N)} = \Phi^{-1}(G_j(r_j + m_j)) \quad j = 1, \dots, J,$$

where  $G_j(\cdot)$  is the marginal distribution function of  $j$ -th element of  $G(\cdot)$ ,  $m_j = \int \epsilon_j dG_j(\epsilon_j)$ , and  $\Phi^{-1}(\cdot)$  denotes the inverse cumulative distribution function of the standard normal distribution. From Theorem 2, we can immediately derive that, if the assumed model agrees with the true model (i.e  $\beta = \tilde{\beta}, \mathbf{V} = \tilde{\mathbf{V}}, G = \tilde{G}$ ), every element  $r_j^{(N)}$  of our normalized surrogate residual will follow a standard normal distribution.

The normalized surrogate residual is similar to the randomized quantile residual proposed by Dunn and Smyth (1996) who transformed the original residual to follow a standard normal distribution under the null hypothesis in generalized linear models with a continuous, binary or count response variable. In fact, for DCMs with a binary outcome where the choice set contains 2 alternatives, our proposed normalized surrogate residual is the same as the randomized quantile residual.

As a numerical illustration, Figure S7 in the supplementary material displays the residual-by-covariates and QQ-plots of our normalized residuals for the improved model of Example 2, and we can see the residual symmetrically distributed around 0, in comparison to the slightly positive-skewed patterns in Figure S1.

### 3.4.2 Bootstrapping Residuals

The patterns of our surrogate residuals contain two sources of uncertainties: (1) modeling error and (2) sampling. The latter is introduced by conditional sampling from  $\mathcal{S}$ . The sampling uncertainty is even more important when the sample size is small. Therefore, we also recommend bootstrapping  $K$  copies of the empirical distributions of the residuals to reduce the sampling uncertainty for a small sample size. For more details about this bootstrapping approach, see Section 5 in Liu and Zhang (2018).

For comparison, We provide the residual-by-covariate plots of the surrogate residual ( $\mathbf{r}$ ) and bootstrapping residual with  $K = 10$  ( $\mathbf{r}^{(B)}$ ) in the setting of Example 2 when the model is specified correctly, with a small sample size of 50. Specifically, three replications of  $\mathbf{r}$  and  $\mathbf{r}^{(B)}$  are realized, and each replication of  $r_1$ -by- $x$  and  $r_1^{(B)}$ -by- $x$  plots is presented in Supplementary Material Figure S8. We omit the second and third dimensions of the residuals because their patterns are similar. It is shown that every replication of the surrogate residual exhibits different shapes. In particular, the second and third replications provide a wave-like shape, which can be a false alarm of a misspecified mean structure. However, all the replications of the bootstrapping residuals are randomly distributed around 0, exhibiting no abnormal patterns.

## 3.5 Simulation Results under Small Sample Scenarios

In this subsection, a Monte Carlo simulation study with three scenarios were performed to investigate the finite-sample properties of Pearson, deviance and surrogate residuals. All three scenarios consider a multinomial logistic regression with a 3-alternative response  $Y = \{1, 2, 3\}$ , with the following utility levels

$$U_{ij} = \beta_{j0} + \beta_{j1}x_i + \epsilon_{ij}, \quad j = 1, 2, 3; i = 1, \dots, n,$$

where  $x_i$  followed  $Unif(0, 1)$ , and  $\beta_{30}$  and  $\beta_{31}$  were fixed to zero as the reference. In Scenario I,  $\beta_{10} = 0.8$ ,  $\beta_{11} = -2$ ,  $\beta_{21} = -1.2$  and  $\beta_{22} = 2$ , which resulted in the empirical probabilities of  $Y = 1, 2, 3$  close to 33%, 33% and 34%, respectively. In Scenario II,  $\beta_{10} = 1.6$ ,  $\beta_{11} = -1.5$ ,  $\beta_{21} = 2.5$  and  $\beta_{22} = -5.5$ , which led to the empirical probabilities of

*Y* = 1, 2, 3 close to 50%, 25% and 25%, respectively. In Scenario III,  $\beta_{10} = 1.1$ ,  $\beta_{11} = 1.1$ ,  $\beta_{21} = -2.2$  and  $\beta_{22} = 4$ , which yielded the empirical probabilities of  $Y = 1, 2, 3$  close to 70%, 15% and 15%, respectively. All results were based on 10,000 Monte Carlo replications and  $n = 20$ . For each scenario, the covariate  $x_i$ 's ( $i = 1, \dots, 20$ ) were held as constants over all replications. In each of the 10,000 replications, we fit the correct model and compute the Pearson's residual  $r_{ij}^{(P)}$ , deviance residual  $r_{ij}^{(D)}$  and surrogate residual  $r_{ij}^{(S)}$ , for  $i = 1, \dots, 20$  and  $j = 1, 2, 3$ . Because the baseline for  $r_{ij}^{(P)}$  and  $r_{ij}^{(D)}$  were standard normal, but the null distribution of  $r_{ij}^{(S)}$  was Gumbel, for easy comparison, we transformed  $r_{ij}^{(S)}$  to the normalized surrogate residual, denoted by  $r_{ij}^{(N)}$ . Therefore, all  $r_{ij}^{(P)}$ ,  $r_{ij}^{(D)}$ ,  $r_{ij}^{(N)}$  had a standard normal baseline.

Tables S1, S2 and S3 in the supplementary material present the empirical mean, variance, skewness and kurtosis coefficients for  $r_{ij}^{(P)}$ ,  $r_{ij}^{(D)}$ ,  $r_{ij}^{(N)}$  in Scenario I, II, III, respectively. If the residual approximately follows the standard normal distribution, the average mean, variance, skewness and kurtosis coefficients should be close to 0, 1, 0 and 3, respectively. In all three scenarios, we can see all four characteristics of the normalized surrogate residual are close to their theoretical values, suggesting the normality of the residuals. The mean and variance of the Pearson's residual are also close to 0 and 1, respectively, but its skewness and kurtosis coefficients are far away from their limiting value, indicating that it may not be well approximated by the standard normal distribution. The deviance residual performed the worst compared with other residuals, as its characteristics were far away from its theoretical values. Tables S1 to S3 also provide the value of the Kolmogorov-Smirnov test (K-S test) statistic and p-value used to test whether the residual follows the standard normal distribution. Under the 5% significance level, the Pearson's and deviance residuals in all three scenarios were rejected from the standard normality hypothesis. However, most of K-S test p-values from the normalized surrogate residuals is greater than 5%, indicating the reasonableness of standard normality.

## 4 Residual for General Models with a Categorical Outcome

In this section, we show that surrogate residuals can be expanded to define residuals for general statistical models with a categorical outcome. Suppose that the assumed model for a categorical outcome  $Y \in \{1, \dots, J\}$  is

$$\Pr(Y = y) = F_a(y; \mathbf{X}, \boldsymbol{\beta}) \quad y = 1, 2, \dots, J, \quad (7)$$

where  $\mathbf{X}$  and  $\boldsymbol{\beta}$  are covariates and unknown parameters, and  $F_a(\cdot)$  is a nominal distribution function with  $J$  categories satisfying  $F_a(j; \mathbf{X}, \boldsymbol{\beta}) \geq 0$  ( $\forall j \in \{1, \dots, J\}$ ) and  $\sum_{j=1}^J F_a(j; \mathbf{X}, \boldsymbol{\beta}) = 1$ . This model is broad enough to cover nearly all regression models with categorical outcomes. For model (7), we can define the hypothetical variable  $\mathbf{Z} = [Z_1, \dots, Z_J]^T$  as

$$Z_j = \ln(F_a(j; \mathbf{X}, \boldsymbol{\beta})) + \epsilon_j \quad j = 1, 2, \dots, J,$$

where  $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_J]^T$  denotes a  $J$ -dimensional noise variable following a multivariate Gumbel distribution with c.d.f  $MG(\boldsymbol{\epsilon}) = \exp(-\sum_{j=1}^J e^{-\epsilon_j})$ . It is evident that  $\epsilon_j$ 's are mutually independent variables following a univariate standard Gumbel distribution. For simplicity, we abbreviate  $\ln(F_a(j; \mathbf{X}, \boldsymbol{\beta}))$  as  $\mu_j$  henceforth. Now, we need to build up a joint distribution of  $Y$  and  $\mathbf{Z}$ . Specifically, after derivation (see Appendix C in the supplementary material for more detail), we can show that

$$\Pr(Z_y > Z_j, \forall j \neq y) \equiv F_a(y; \mathbf{X}, \boldsymbol{\beta}) \quad \forall y = 1, \dots, J.$$

Thus the joint distribution can be determined by setting  $Y \triangleq y$  if and only if  $Z_y > Z_j (\forall j \neq y)$ . Now, let the surrogate variable  $\mathbf{S}$  be the conditional distribution of  $\mathbf{Z}$  given  $Y$ . More specifically,  $\mathbf{S} \sim \mathbf{Z}|(Y, \mathbf{X})$  follows a truncated distribution obtained by truncating the distribution of  $\mathbf{Z}|\mathbf{X} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$  using the region  $Z_Y > Z_j (\forall j \neq Y)$ , where  $\boldsymbol{\mu} = [\mu_1, \dots, \mu_J]^T$ .

Finally, we can define the surrogate residual  $\mathbf{R}$  as

$$\mathbf{R} = \mathbf{S} - E_0(\mathbf{S}|\mathbf{X}) = \mathbf{S} - E(\mathbf{Z}|\mathbf{X}) = \mathbf{S} - \boldsymbol{\mu} - \int \boldsymbol{\epsilon} dMG(\boldsymbol{\epsilon}), \quad (8)$$

where  $\int \epsilon dMG(\epsilon) \approx 0.5772$  is the expectation of  $MG(\cdot)$ . In practice, given the observed data  $\{\mathbf{x}, y\}$  and a fitted model  $F_a(y; \mathbf{X}, \hat{\beta})$ , we draw a sample  $\mathbf{s}$  from  $\mathbf{S}$  by plugging in the estimated parameter  $\hat{\beta}$ , then let  $\mathbf{r} = \mathbf{s} - \hat{\mu} - 0.5772$  be the surrogate residual. If the assumed model agrees with the true model, say the true model for  $Y$  is  $\tilde{F}(y; \mathbf{X}, \tilde{\beta})$ , the theorem below summarizes the properties of  $\mathbf{R}$ .

**Theorem 5.** *If the assumed model (7) agrees with the underlying true model (i.e.,  $F_a = \tilde{F}$ ,  $\beta = \tilde{\beta}$ ), the following results hold*

(1) *The surrogate variable  $\mathbf{S}$  has the same distribution as  $\mathbf{Z}$ , i.e.  $\mathbf{S} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\epsilon}$  follows the distribution  $MG(\boldsymbol{\epsilon})$ .*

(2) *The residual variable  $\mathbf{R}$ , independent from  $\mathbf{X}$ , follows the distribution  $MG(\mathbf{r} + \mathbf{m})$ , i.e.  $Pr(\mathbf{R} \leq \mathbf{r}) = Pr(\mathbf{R} \leq \mathbf{r} | \mathbf{X}) = MG(\mathbf{r} + \mathbf{m})$ , where  $\mathbf{m} = \int \epsilon dMG(\epsilon) \approx 0.5772$ . As a result,  $E(\mathbf{R} | \mathbf{X}) = \mathbf{0}$ , and  $Var(\mathbf{R} | \mathbf{X}) = Var(\boldsymbol{\epsilon}) \approx 1.6449 \times I_{J \times J}$ , where  $I_{J \times J}$  denotes a  $J \times J$  identity matrix.*

Proof of Theorem 5 is similar to the one discussed in Section 3.2. Theorem 5(2) shows that the residual defined in (8) is similar to the residual defined for DCMs, has zero mean, homogeneous variance, and an explicit reference distribution. Therefore, it can be used for model diagnostics in a similar way to the surrogate residual defined for DCMs.

**Remark 3.** *For multinomial logistic regression, the residual defined in (8) is exactly equal to the residual given in Section 3.1. This reveals an interesting connection of the surrogate residual for general models with a categorical outcome to that for DCMs.*

## 5 MPN Data Analysis

We present a real data analysis to illustrate how to perform residual diagnosis in practice using the Netherlands Mobility Panel (MPN) data. The data were collected to investigate the travel behavior of Dutch residents. Since the MPN is a mobility panel, multiple waves are available since 2013. As an illustrative example, we only make use of the data from the first wave, which includes a total of 35,257 travel records from 5,297 individuals. See Hoogendoorn-Lanser et al. (2015) for more details.

1  
2  
3 For travel behavior analysis, we focus on a categorical outcome containing six different  
4 travel modes (car as driver, car as passenger, bicycle, ebike, public transport, and walking).  
5 We use the discrete choice model to analyze the travel modes, and focus on the following  
6 covariates: travel distance (kilometers), departure time (day time, night time), driving  
7 license (with, without), gender, age, education (low, medium, high), work situation, and  
8 household income (low, medium, rich). All the covariates are individual-specific variables.  
9 Now, we illustrate how to use our surrogate residual to diagnose and improve model fitting.  
10  
11

12 Used in our initial analysis is a multinomial logistic regression, including all above-  
13 mentioned covariates in linear terms. We first investigate the association between the travel  
14 distance and travel mode choice, since the distance from the destination exerts direct and  
15 may be the most important influence to people's decision on which kind of transport to  
16 choose. Figure S9 in the supplementary material plots residual versus travel distance (we  
17 take travel distance under the logarithmic coordinate for better display effect). We see that  
18 the mean of scatter points in plots (a), (e) and (f), corresponding to the utilities of car as  
19 driver, public transport and walking, show a wave shape. Especially the points in Figure  
20 S9(a) have a positive mean shift to the right of the vertical dashed line; on the contrary,  
21 the points in Figure S9(e) have a negative mean shift in the same interval. Plus, plot  
22 (c) of Figure S9 reveals a wave shape, implying a misspecified mean structure. All these  
23 characteristics indicate that a linear form of travel distance is inadequate. We therefore  
24 tried using a natural logarithm form, quadratic form and piecewise linear form of travel  
25 distance in our multinomial logistic model, and found that the natural logarithm form  
26 could most effectively alleviate the mean shift. Therefore, we replaced the linear form of  
27 distance with a natural logarithm form from subsequent analysis. The updated residual-  
28 by-log(distance) plots are shown in Figure S10, which effectively smooth out the mean  
29 shifts. In fact, our updated model improved McFadden R-squared from 0.2581 to 0.3067.  
30 Moreover, a likelihood ratio test with a p-value less than  $2.2 \times 10^{-16}$  also suggests that a  
31 logarithm term of distance resulted in a better fit of the data.  
32  
33

34 Our further examination reveals that there exists an interaction effect between the  
35 driving license and travel distance. Specifically, we fit two models using the updated mean  
36

1  
2  
3 structure based on two groups. The first group is the subjects who have a driving license,  
4 while the second group is the subjects who do not have a driving license. We draw the  
5 plots of surrogate response versus log(distance) for above two groups in Figure S11. If  
6 the true model does not contain an interaction term  $\log(\text{distance}) \times \text{license}$ , the slopes of  
7 the surrogate responses by  $\log(\text{distance})$  should be all the same between these two data  
8 groups. However, our plots of surrogate responses by  $\log(\text{distance})$  present different slopes.  
9 Therefore, we introduce a new interaction term  $\log(\text{distance}) \times \text{license}$  to our model. This  
10 adjustment slightly improves McFadden R-squared from 0.3067 to 0.3073. The likelihood  
11 ratio test for including this interaction effect compared to the former model yields a p-value  
12 of  $2.0 \times 10^{-12}$ .  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

## 23 6 Discussion

## 24

25 In this article, we proposed the surrogate variables and residuals for the discrete choice  
26 models. We defined the surrogate residual, and comprehensively analyzed its theoretical  
27 and graphical properties. Numerical examples in a variety of settings confirm that our  
28 residual can be informative and useful in detecting misspecified mean structure, checking  
29 individual-specific coefficients, and detecting interaction effects. Our surrogate approach  
30 overcomes the problems in the Pearson's, deviance and other nonrandomized residuals, and  
31 enriches the diagnostic methods in DCMs. We extend our residuals for more general models  
32 with an unordered categorical response. Finally, our analysis of the MPN data confirms  
33 that our surrogate residual can help us improve the model fitting in practice.  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44

## 45 Acknowledgments

## 46

47 Zhang's research was partially supported by grants NIH R01 MH116527 from National  
48 Institute of Mental Health and NSF DMS-1722544 from the National Science Foundation  
49 of United States of America.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Supplementary Material

**Appendices and supplementary figures:** “residual\_DCM\_supp.pdf” presents appendices (Appendices A, B and C) and supplementary tables and figures unshown in the manuscript.

**R-code for numerical examples:** The R code for the numerical examples are shown at “residual\_DCM\_codes.R”.

**MPN data:** The real data set used in Section 5 is available at <https://www.mpnldata.nl>.

## References

- Ben-Akiva, M. E. and S. R. Lerman (1985). *Discrete choice analysis: theory and application to travel demand*, Volume 9. MIT press.
- Chan, Y. H. (2005). Biostatistics 305. multinomial logistic regression. *Singapore Medical Journal* 46(6), 259–269.
- Dunn, P. K. and G. K. Smyth (1996). Randomized quantile residuals. *Journal of Computational and Graphical Statistics* 5(3), 236–244.
- Goeman, J. J. and S. le Cessie (2006). A goodness-of-fit test for multinomial logistic regression. *Biometrics* 62(4), 980–985.
- Gupta, A. K., T. Nguyen, and L. Pardo (2008). Residuals for polytomous logistic regression models based on  $\varphi$ -divergences test statistics. *Statistics* 42(6), 495–514.
- Hausman, J. A. and D. McFadden (1981). Specification tests for the multinomial logit model. *Econometrica: Journal of the Econometric Society*, 1219–1240.
- Hoogendoorn-Lanser, S., N. T. Schaap, and M. J. OldeKalter (2015). The netherlands mobility panel: An innovative design approach for web-based longitudinal travel data collection. *Transportation Research Procedia* 11, 311–329.

- 1  
2  
3 Koppelman, F. S. and C. Bhat (2006). A self instructing course in mod-  
4  
e choice modeling: multinomial and nested logit models. OnlineDoc:  
5  
http://www.caee.utexas.edu/prof/bhat.  
6  
7  
8 Liu, D. G. and H. P. Zhang (2018). Residuals and diagnostics for ordinal regression models:  
9  
A surrogate approach. *Journal of the American Statistical Association* 113(522), 845–  
10  
854.  
11  
12 Lovreglio, R., D. Borri, D. Luigi, and A. Ibeas (2014). A discrete choice model based on  
13  
random utilities for exit choice in emergency evacuations. *Safety Science* 62, 418–426.  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
59  
60  
Nagel, H. and R. Hatzinger (1992). Diagnostics in some discrete choice models. *Computa-  
tional Statistics & Data Analysis* 13(4), 407–414.  
Pagan, A. and F. Vella (1989). Diagnostic tests for models based on individual data: a  
survey. *Journal of Applied Econometrics* 4(S1), S29–S59.  
Train, K. E. (1998). Recreation demand models with taste differences over people. *Land  
Economics*, 230–239.

# Supplementary Material: Surrogate Residuals for Discrete Choice Models

Chao Cheng

Department of Mathematical Sciences, Tsinghua University  
Rui Wang

Department of Public Economics, School of Economics, Xiamen University  
and  
Heping Zhang\*

Department of Biostatistics, School of Public Health, Yale University

August 26, 2019

This file presents the appendices (Appendices A, B and C) and supplementary tables and figures unshown in the manuscript.

## Appendix A: Generating Samples from Truncated Multi-dimensional Distribution

In this section, we introduce a Gibbs method to generate samples from truncated multi-dimensional distribution like  $\mathbf{S}$  in section 3.1. Before the multi-dimensional case, we at first present the method to generate truncated univariate distribution, because the ability to generate univariates is a central building block in the Gibbs method.

### 1. The univariate case

Let  $X$  follow a distribution  $G$  truncating on the interval  $(a, b)$ , where  $G$  is a univariate continuous cumulative distribution function,  $a$  is the left truncated point and  $b$  is the right truncated

\*Address for correspondence: Heping Zhang, Department of Biostatistics, School of Public Health, Yale University, New Haven, CT, USA. Email: heping.zhang@yale.edu

point. The two truncated points are arbitrary but fixed on  $\mathbb{R}$  and  $-\infty \leq a < b \leq +\infty$ . Then, the distribution function of  $X$  is

$$F(x) = \frac{G(x) - G(a)}{G(b) - G(a)} \quad a \leq x \leq b.$$

Based on the process of *inverse CDF method*, solving the equation  $F(x) = u$ ,  $u \in (0, 1)$ , we get

$$x = G^{-1}\left(u(G(b) - G(a)) + G(a)\right).$$

Then, we simulate  $X$  as follows:

**Step 1:** Generate a sample  $u$  from  $\text{Unif}(0, 1)$ , where  $\text{Unif}(0, 1)$  denotes the standard uniform distribution.

**Step 2:** Return  $X = G^{-1}\left(u(G(b) - G(a)) + G(a)\right)$ . Then,  $X$  follows the target distribution.

## 2. The multivariate case

Now, we propose a Gibbs sampling approach for drawing samples from a truncated multi-dimensional density like  $\mathbf{S}$  in section 3.1. Gibbs sampling is a Markov chain Monte Carlo (MCMC) method for generation of samples from high dimensional densities by drawing the samples from full conditional densities. In other words, to implement the Gibbs sampling for an  $J$ -dimensional vector, we need to find full conditionals of its subvectors.

Let  $\mathbf{X} = [X_1, \dots, X_J]^T$  denote a variable that follows a  $J$ -dimensional distribution  $G(x_1, \dots, x_J)$  using the support area  $\mathcal{D}_k$ , where  $G(x_1, \dots, x_J)$  is a continuous cumulative distribution function and  $\mathcal{D}_k$  is a support area in  $\mathbb{R}^J$  such as  $X_k > X_j$ ,  $\forall j \neq k$ . Denote by  $\mathbf{x} = (x_1, x_2, \dots, x_J)^T$  a sample from  $\mathbf{X}$ . We shall denote the c.d.f of the full conditionals  $X_j|X_{-j}$ ,  $j = 1, 2, \dots, J$ , where  $X_{-j} = (X_1, X_2, \dots, X_{j-1}, X_{j+1}, \dots, X_J)^T$ , by  $G_{X_j|X_{-j}}(x_j)$ . The conditional density  $X_j|X_{-j}$  is

$$F_{X_j|X_{-j}}(x_j|x_{-j}) = \frac{G_{X_j|X_{-j}}(x_j|x_{-j}) - G_{X_j|X_{-j}}(a|x_{-j})}{G_{X_j|X_{-j}}(b|x_{-j}) - G_{X_j|X_{-j}}(a|x_{-j})},$$

1 where  
 2  
 3  
 4  
 5

$$\begin{cases} a = \max x_{-j}, b = +\infty & \text{if } j = k \\ a = -\infty, b = x_k & \text{if } j \neq k \end{cases}$$

6 and  
 7  
 8  
 9

$$G_{X_j|X_{-j}}(x_j|x_{-j}) = \int_{-\infty}^{x_j} \frac{g(x_1, \dots, x_{j-1}, v_j, x_{j+1}, \dots, x_J)}{\int_{-\infty}^{+\infty} g(x_1, \dots, x_{j-1}, v_j, x_{j+1}, \dots, x_J) dv_j} dv_j.$$

10 Here,  $g(\cdot)$  is the corresponding p.d.f of  $G(\cdot)$ . The distribution above is nothing but a truncated  
 11 univariate distribution whose samples can be generated by the previous subsection. Now we present  
 12 the Gibbs sampling scheme that will be used to produce samples from  $\mathbf{X}$ . First, we get initial  
 13 values  $\{x_k^{(0)}\}_{k=2}^J$  from the support area  $\mathcal{D}_k$ , and then for  $i = 1, 2, \dots, M$  generate samples according  
 14 to  
 15

$$\begin{aligned} X_1|x_2^{(i-1)}, x_3^{(i-1)}, \dots, x_J^{(i-1)} &\longrightarrow x_1^{(i)} \\ X_2|x_1^{(i)}, x_3^{(i-1)}, \dots, x_J^{(i-1)} &\longrightarrow x_2^{(i)} \\ &\dots \\ X_{J-1}|x_1^{(i)}, \dots, x_{J-2}^{(i)}, x_J^{(i-1)} &\longrightarrow x_{J-1}^{(i)} \\ X_J|x_1^{(i)}, x_2^{(i)}, \dots, x_{J-1}^{(i)} &\longrightarrow x_J^{(i)} \end{aligned}$$

26 where the arrows represent sampling from the corresponding distribution, and each  $x_k$  is sampled  
 27 according to the univariate scheme. We should discard the first  $B$  samples, where  $B$  is a given  
 28 number in advance, which are called burn-in samples, and the remaining samples are considered  
 29 as the samples from  $\mathbf{X}$ .  
 30  
 31  
 32  
 33  
 34  
 35  
 36  
 37

## 38 Appendix B: Proofs for Theorems 3 and 4

41 In this section, we prove Theorems 3 and 4. First, we derive the distribution of  $\hat{\mathbf{S}}$  that follows  
 42 the truncated distribution obtained by truncating the distribution of  $\hat{\mathbf{U}} = \hat{\mathbf{V}}(\mathbf{X}, \hat{\boldsymbol{\beta}}) + \epsilon$  using the  
 43 area  $\{\hat{U}_y > \hat{U}_j, \forall j \neq y\}$ , denoted by  $\hat{\mathcal{D}}_y$  henceforth, given  $Y = y$ . Noting that  $\{\hat{U}_y > \hat{U}_j, \forall j \neq y\}$   
 44 only represents the order information, we have  $\hat{\mathcal{D}}_y \equiv \mathcal{D}_y$ , whether  $\hat{\boldsymbol{\beta}}$  is a consistent estimate or not.  
 45 When  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + o_p(1)$  and under the assumption that  $\mathbf{V}(\mathbf{X}, \boldsymbol{\beta})$  is a continuous function, we have  
 46  $\mathbf{V}(\mathbf{X}, \hat{\boldsymbol{\beta}}) = \mathbf{V}(\mathbf{X}, \boldsymbol{\beta}) + o_p(1)$ , thus  $\hat{\mathbf{U}} = \mathbf{U} + o_p(1)$ . Therefore, for all continuity sets  $\mathcal{B} \in \mathbb{R}^J$ , we  
 47 have  $Pr(\hat{\mathbf{U}} \in \mathcal{B}) = Pr(\mathbf{U} \in \mathcal{B}) + o(1)$ , which implies  $Pr(\hat{\mathbf{U}} \in \mathcal{C} \cap \mathcal{D}_y) = Pr(\mathbf{U} \in \mathcal{C} \cap \mathcal{D}_y) + o(1)$

and  $\Pr(\hat{\mathbf{U}} \in \mathcal{D}_y) = \Pr(\mathbf{U} \in \mathcal{D}_y) + o(1)$ . Now, conditional on  $\mathbf{X}$ , for any arbitrary but fixed  $\mathbf{c} = [c_1, \dots, c_J]$ ,

$$\begin{aligned}
& |Pr(\hat{\mathbf{S}} \leq \mathbf{c}) - Pr(\mathbf{S} \leq \mathbf{c})| \\
&= \left| \sum_{y=1}^J Pr(\hat{\mathbf{S}} \leq \mathbf{c}|Y=y) Pr(Y=y) - \sum_{y=1}^J Pr(\mathbf{S} \leq \mathbf{c}|Y=y) Pr(Y=y) \right| \\
&\leq \sum_{y=1}^J \left| Pr(\hat{\mathbf{S}} \leq \mathbf{c}|Y=y) - Pr(\mathbf{S} \leq \mathbf{c}|Y=y) \right| Pr(Y=y) \\
&= \sum_{y=1}^J \left| \frac{Pr(\hat{\mathbf{U}} \in \mathcal{C} \cap \hat{\mathcal{D}}_y)}{Pr(\hat{\mathbf{U}} \in \hat{\mathcal{D}}_y)} - \frac{Pr(\mathbf{U} \in \mathcal{C} \cap \mathcal{D}_y)}{Pr(\mathbf{U} \in \mathcal{D}_y)} \right| Pr(Y=y) \\
&= \sum_{y=1}^J \left| \frac{Pr(\hat{\mathbf{U}} \in \mathcal{C} \cap \mathcal{D}_y)}{Pr(\hat{\mathbf{U}} \in \mathcal{D}_y)} - \frac{Pr(\mathbf{U} \in \mathcal{C} \cap \mathcal{D}_y)}{Pr(\mathbf{U} \in \mathcal{D}_y)} \right| Pr(Y=y) \\
&= \sum_{y=1}^J \left| \frac{Pr(\mathbf{U} \in \mathcal{C} \cap \mathcal{D}_y) + o(1)}{Pr(\mathbf{U} \in \mathcal{D}_y) + o(1)} - \frac{Pr(\mathbf{U} \in \mathcal{C} \cap \mathcal{D}_y)}{Pr(\mathbf{U} \in \mathcal{D}_y)} \right| Pr(Y=y) \\
&= o(1).
\end{aligned}$$

Thus, we can establish  $\Pr(\hat{\mathbf{R}}_{\hat{\beta}} \leq \mathbf{r}|\mathbf{X}) = \Pr(\mathbf{R} \leq \mathbf{r}|\mathbf{X}) + o(1)$ . The unconditional distribution of  $\hat{\mathbf{R}}_{\hat{\beta}}$

$$\Pr(\hat{\mathbf{R}}_{\hat{\beta}} \leq \mathbf{r}) = \int \Pr(\hat{\mathbf{R}}_{\hat{\beta}} \leq \mathbf{r}|\mathbf{x}) d\mu(\mathbf{x}) \rightarrow \int \Pr(\mathbf{R} \leq \mathbf{r}|\mathbf{x}) d\mu(\mathbf{x}) = \Pr(\mathbf{R} \leq \mathbf{r}),$$

based on Lebesgue's Dominated Convergence Theorem. Now, Theorem 3 has been proved. Moreover, based on Theorem 3, we have  $E(\hat{\mathbf{R}}_{\hat{\beta}}|\mathbf{X}) = E(\mathbf{R}|\mathbf{X}) + o(1) = o(1)$  and  $Var(\hat{\mathbf{R}}_{\hat{\beta}}|\mathbf{X}) = Var(\mathbf{R}|\mathbf{X}) + o(1)$ . This completes the proof of Theorem 4.

## Appendix C: Proof for $\Pr(Z_y > Z_j; \forall j \neq y) = F_a(y; \mathbf{X}, \boldsymbol{\beta})$

We now prove that under the assumed model (4.11),  $\Pr(Z_y > Z_j; \forall j \neq y) = F_a(y; \mathbf{X}, \boldsymbol{\beta})$ . At first, abbreviating  $\Pr(Z_y > Z_j; \forall j \neq y)$  as  $p_y$ , we have

$$\begin{aligned}
p_y &= \Pr(Z_y > Z_j, \forall j \neq y) \\
&= \Pr(\mu_y + \epsilon_y > \mu_j + \epsilon_j, \forall j \neq y) \\
&= \Pr(\epsilon_y < \mu_y - \mu_j + \epsilon_j, \forall j \neq y).
\end{aligned}$$

If  $\epsilon_y$  is considered given, this equation is the cumulative density function for each  $\epsilon_j$  evaluated at  $\mu_y - \mu_j + \epsilon_y$ , according to the c.d.f of standard Gumbel distribution, which is  $\exp(-\exp(-(mu_y - mu_j + epsilon_y)))$ . Since all  $\epsilon_y$  are independent, the cumulative distribution over all  $y \neq j$  is the product

1 of the individual cumulative distributions  
 2  
 3  
 4  
 5  
 6

$$p_y|\epsilon_y = \prod_{j \neq y} \exp(-\exp(-(\mu_y - \mu_j + \epsilon_y))).$$

7 But in this question,  $\epsilon_y$  is unknown, so the probability  $p_y$  is the integral of  $p_y|\epsilon_y$  weighted by the  
 8 density of  $\epsilon_y$   
 9  
 10

$$p_y = \int \left( \prod_{j \neq y} e^{-\exp(-(\mu_y - \mu_j + \epsilon_y))} \right) e^{-\epsilon_y - \exp(-\epsilon_y)} d\epsilon_y. \quad (1)$$

11 Mark  $\epsilon_y$  as  $t$  and note that  $\mu_y - \mu_y = 0$ . Thus, collecting the terms in the exponent of  $e$ , we have  
 12  
 13

$$\begin{aligned} p_y &= \int_{t=-\infty}^{+\infty} \left( \prod_j e^{-\exp(-(\mu_y - \mu_j - t))} \right) e^{-t} dt \\ &= \int_{t=-\infty}^{+\infty} \exp\left(-\sum_j e^{-\mu_y + \mu_j + t}\right) e^{-t} dt \\ &= \int_{t=-\infty}^{+\infty} \exp\left(-e^{-t} \sum_j e^{-\mu_y + \mu_j}\right) e^{-t} dt. \end{aligned}$$

24 Define  $s = e^{-t}$  such that  $e^{-t} dt = -ds$ , using the new term, we have  
 25  
 26

$$\begin{aligned} p_y &= \int_0^{+\infty} \exp\left(-s \sum_j e^{-\mu_y + \mu_j}\right) dt = \left( \frac{\exp\left(-s \sum_j e^{-\mu_y + \mu_j}\right)}{-\sum_j e^{-\mu_y + \mu_j}} \right)_0^{+\infty} \\ &= \frac{\exp(\mu_y)}{\sum_{j=1}^J \exp(\mu_j)} = \frac{F_a(y; \mathbf{X}, \boldsymbol{\beta})}{\sum_{j=1}^J F_a(j; \mathbf{X}, \boldsymbol{\beta})} \\ &= F_a(y; \mathbf{X}, \boldsymbol{\beta}). \end{aligned}$$

35 This completes the proof.  
 36  
 37  
 38

## 39 Supplementary Tables and Figures 40







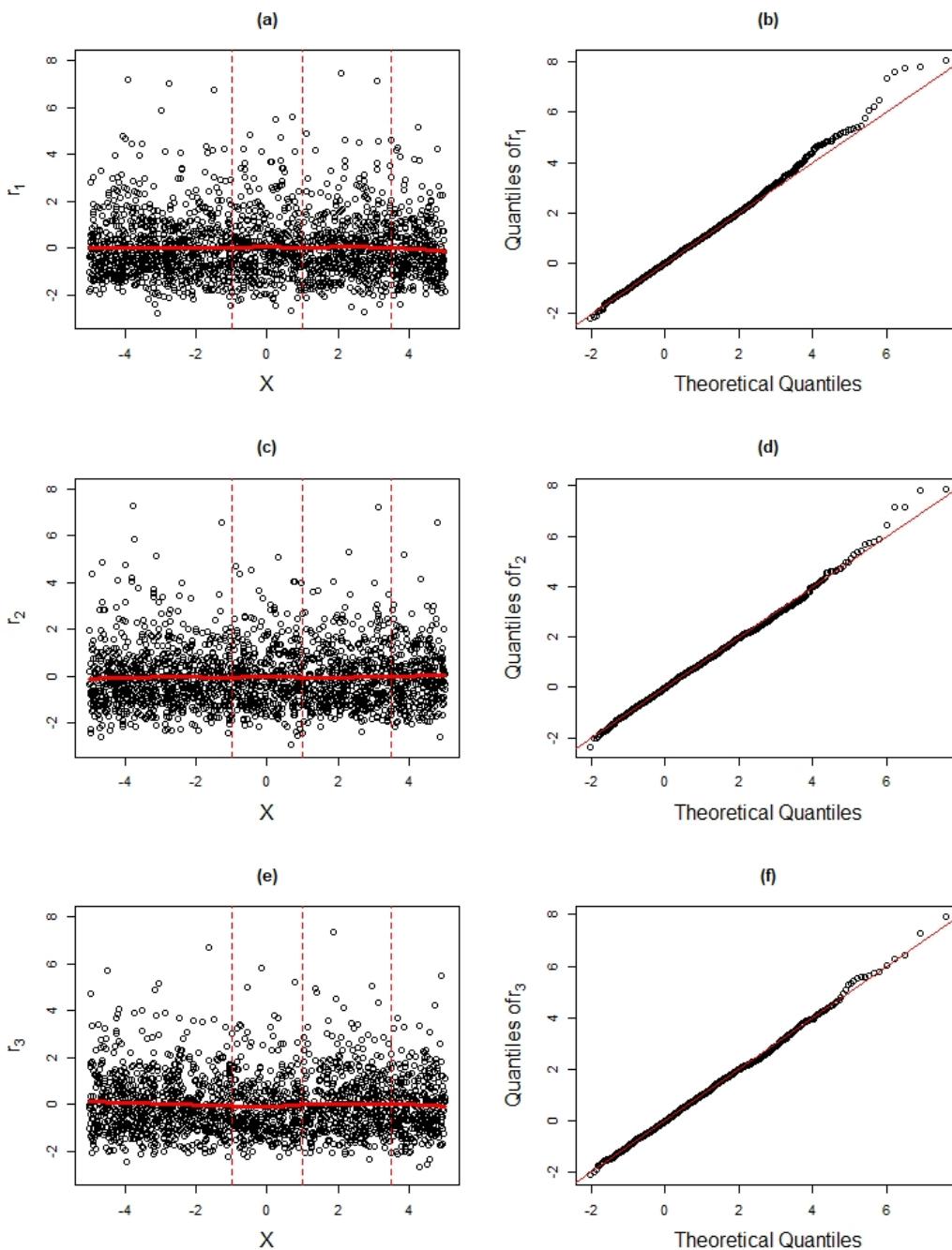


Figure S1: Model diagnostics for the improved model in Example 2 with the surrogate residual. Figures (a), (c) and (e) show the residuals-by-covariate plots of  $r_1$ ,  $r_2$  and  $r_3$  respectively; vertical dashed lines  $X = -1$ ,  $1$ , and  $3.5$ , and loess curves (red solid) are added as references. Figures (b), (d) and (f) are the corresponding QQ-plots (quantile-quantile plots) of above mentioned residuals.

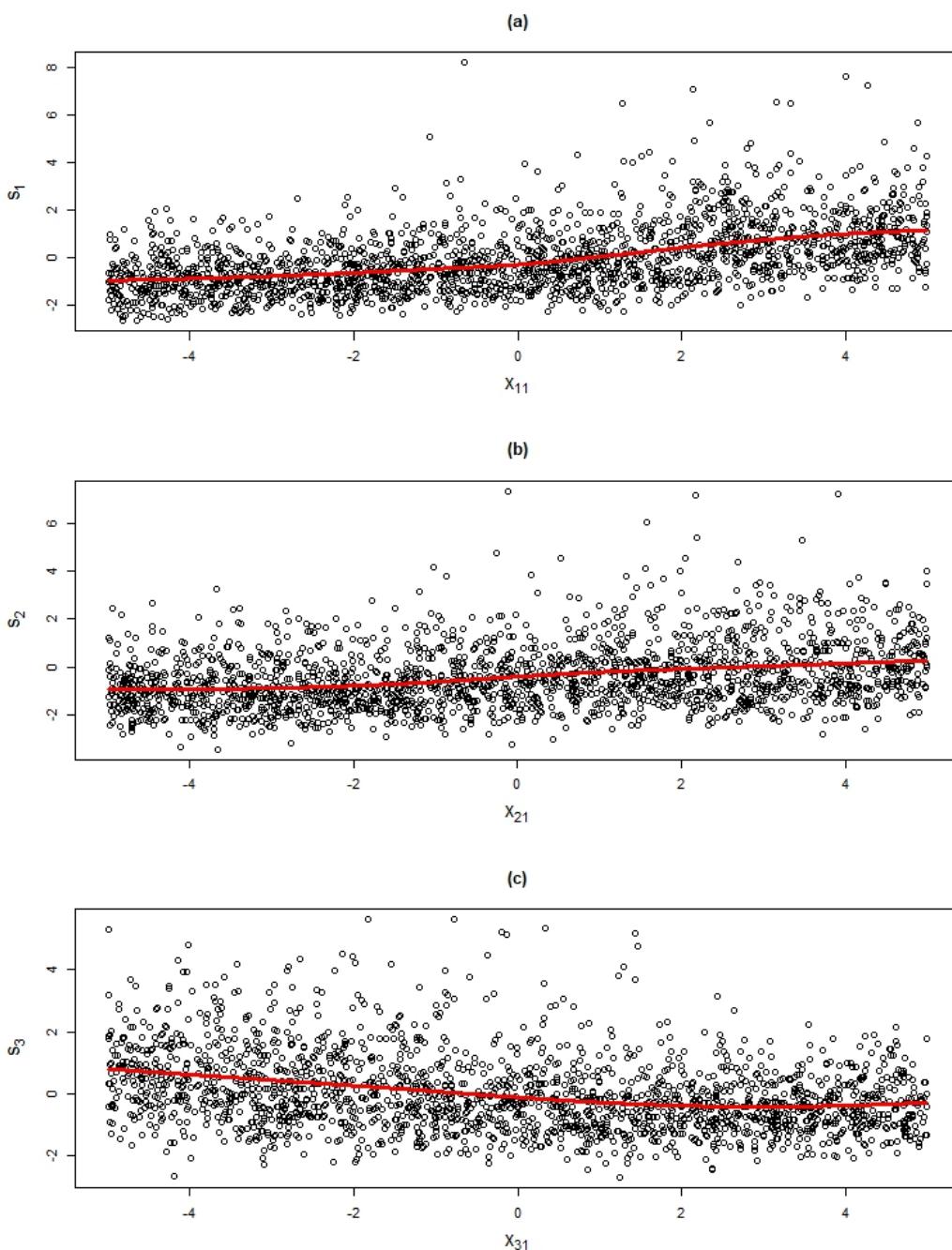


Figure S2: Scatterplots of the surrogate response  $s$  versus  $x_1$  with a nonparametric loess smooth (red line). (a):  $s_1$ -by- $x_{11}$ ; (b):  $s_2$ -by- $x_{21}$ ; (c):  $s_3$ -by- $x_{31}$ .

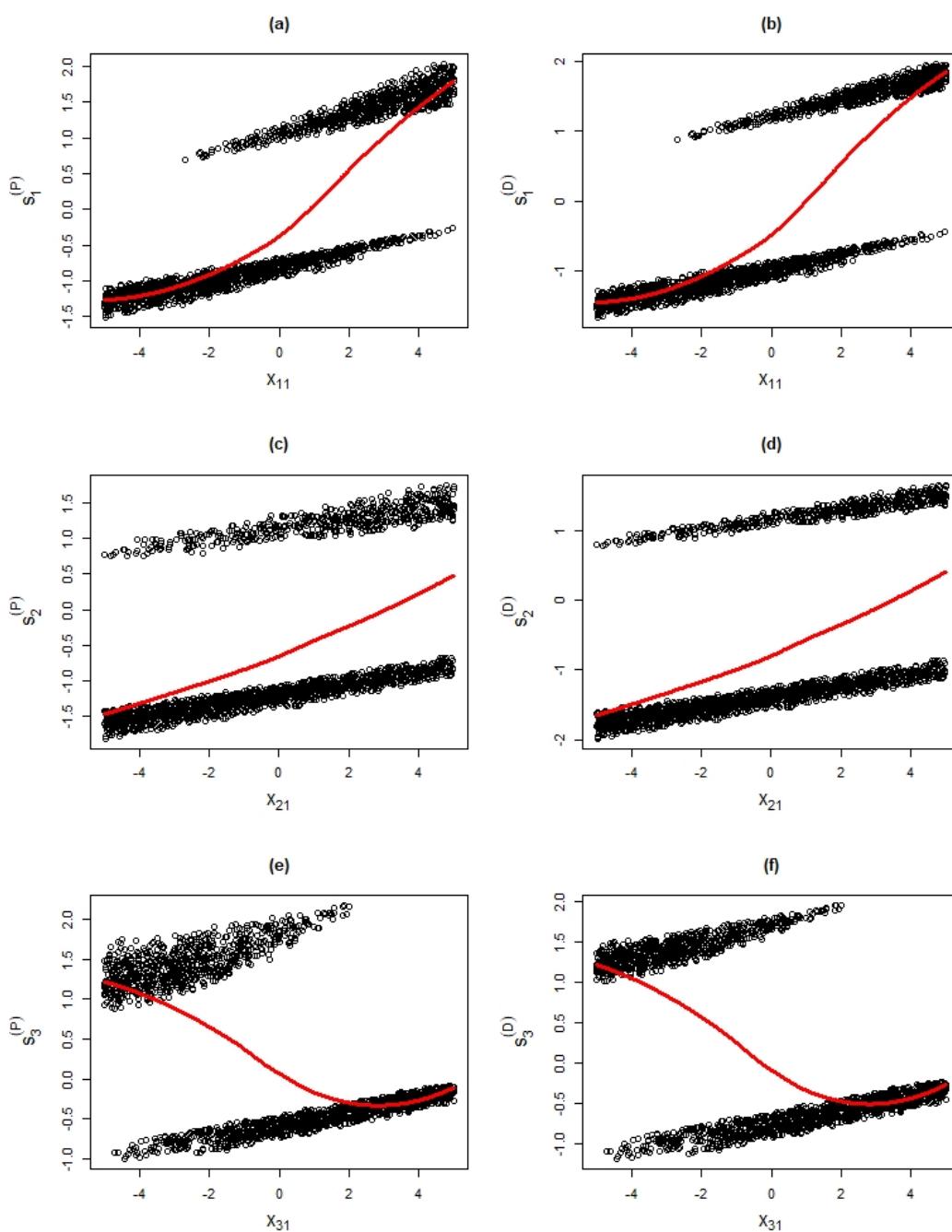


Figure S3: Scatterplots of the surrogate response by Pearson's residual ( $s^{(P)}$ ) and deviance residual ( $s^{(D)}$ ) versus  $x_1$  with a nonparametric loess smooth (red line). Left column:  $s^{(P)}$  by  $x_1$ ; Right column:  $s^{(D)}$  by  $x_1$ .

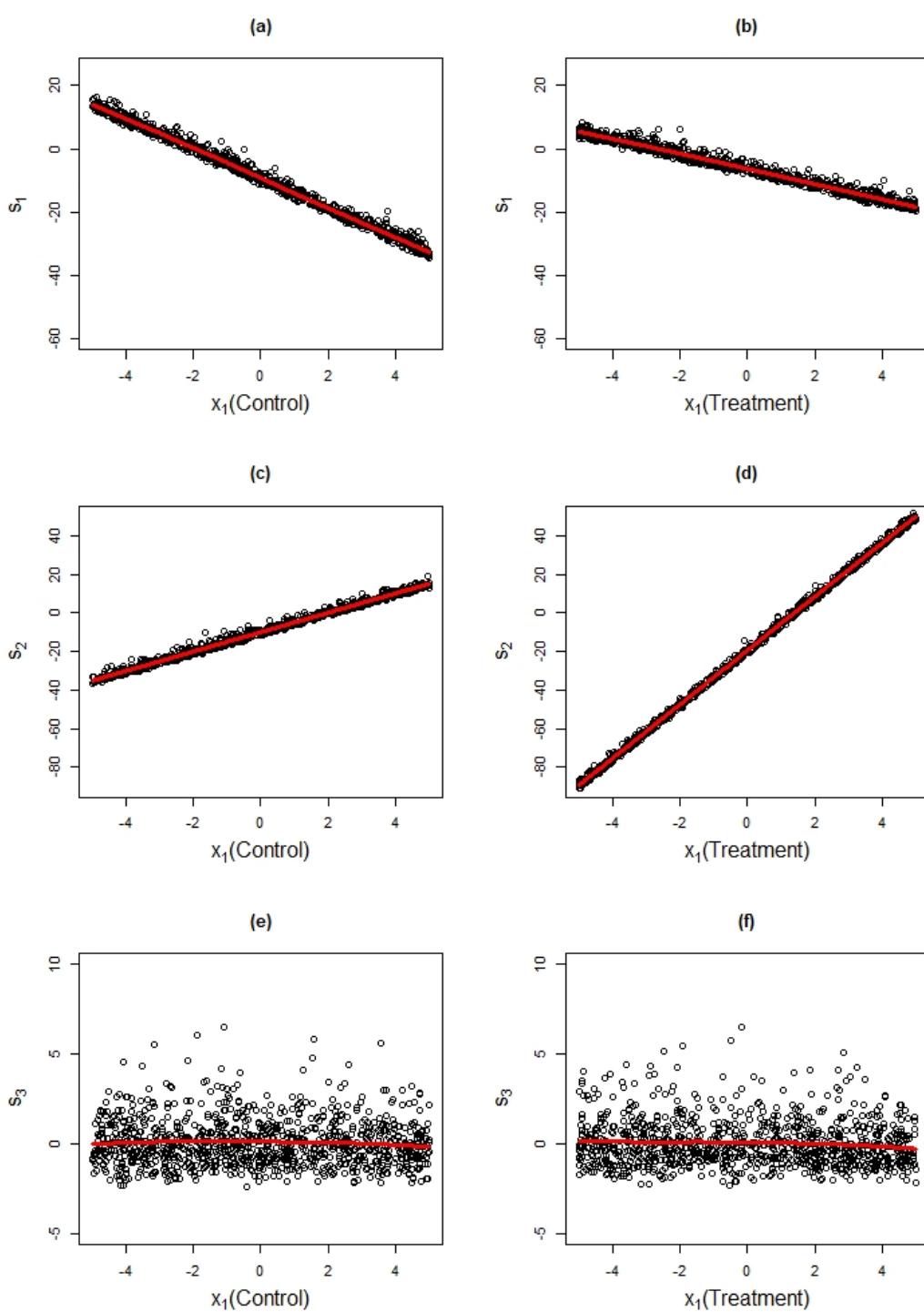


Figure S4: Scatterplots of the surrogate response  $\mathbf{s}$  versus  $x_1$  with a nonparametric loess smooth (red line). Left column: control group; Right column: treatment group.

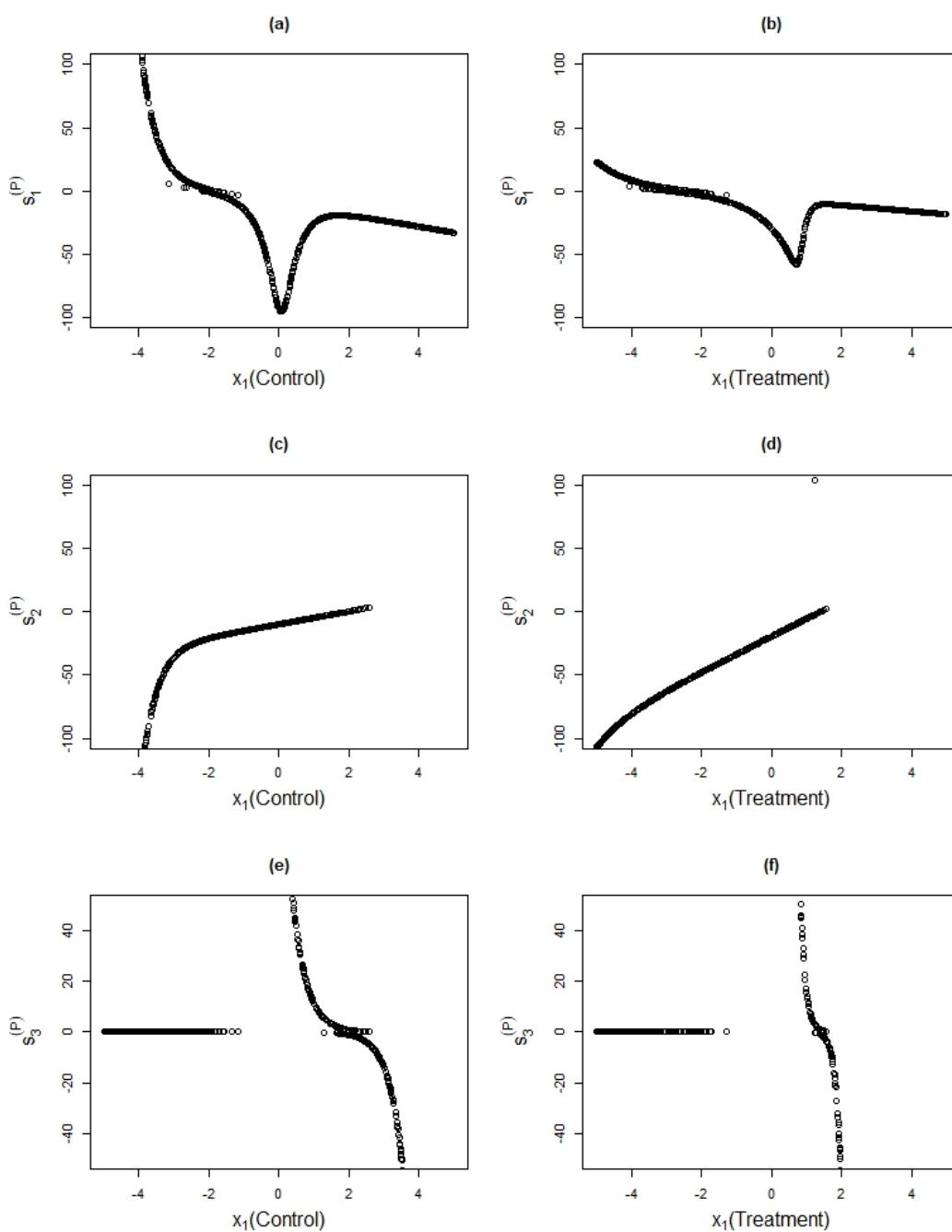


Figure S5: Scatterplots of the surrogate response of Pearson's residual ( $s^{(P)}$ ) versus  $x_1$ . Left column: control group; Right column: treatment group. The loess curve was removed, as some points tend to infinity and a loess curve is hard to compute.

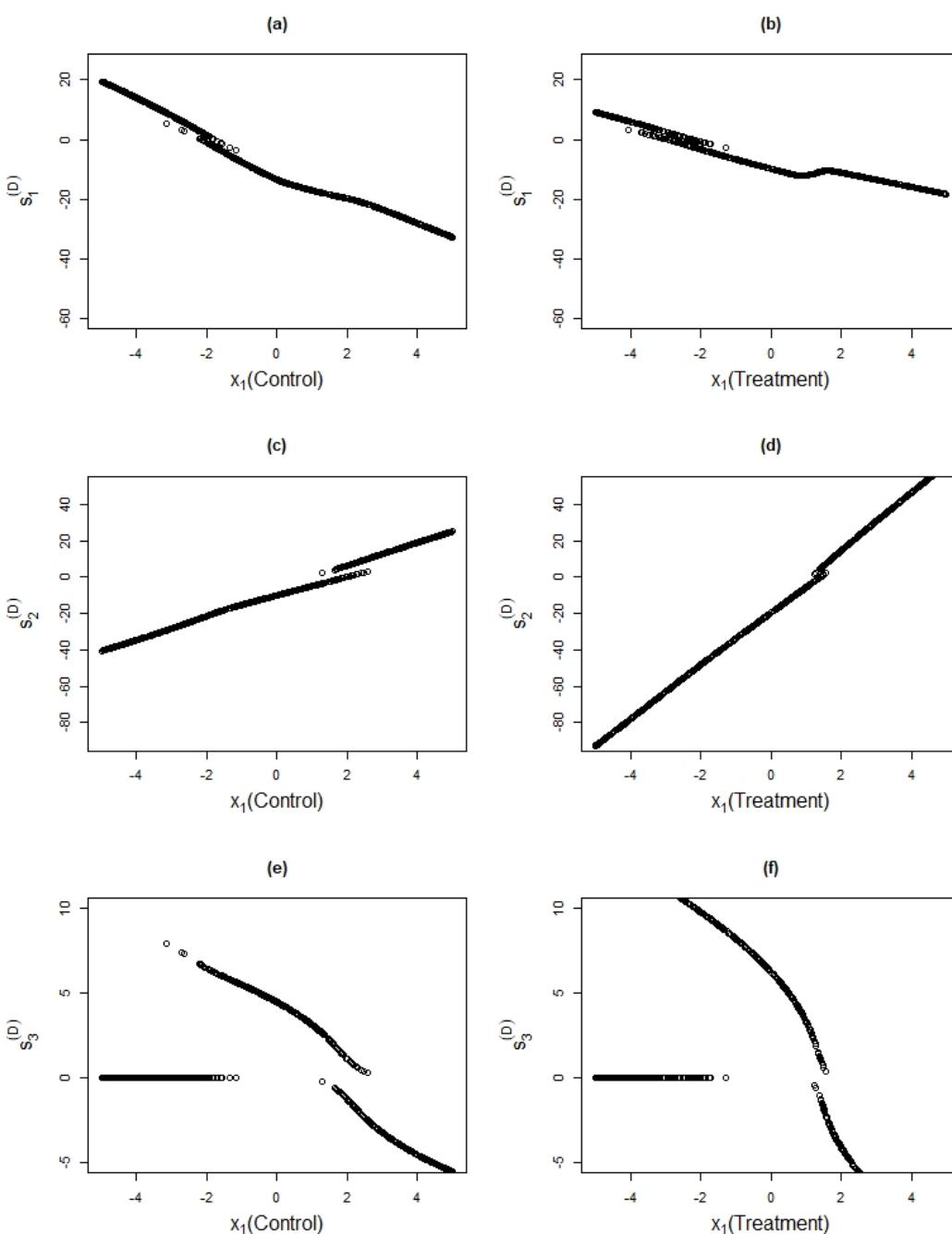


Figure S6: Scatterplots of the surrogate response of deviance residual ( $s^{(P)}$ ) versus  $x_1$ . Left column: control group; Right column: treatment group. The loess curve was removed, as some points tend to infinity and a loess curve is hard to compute.

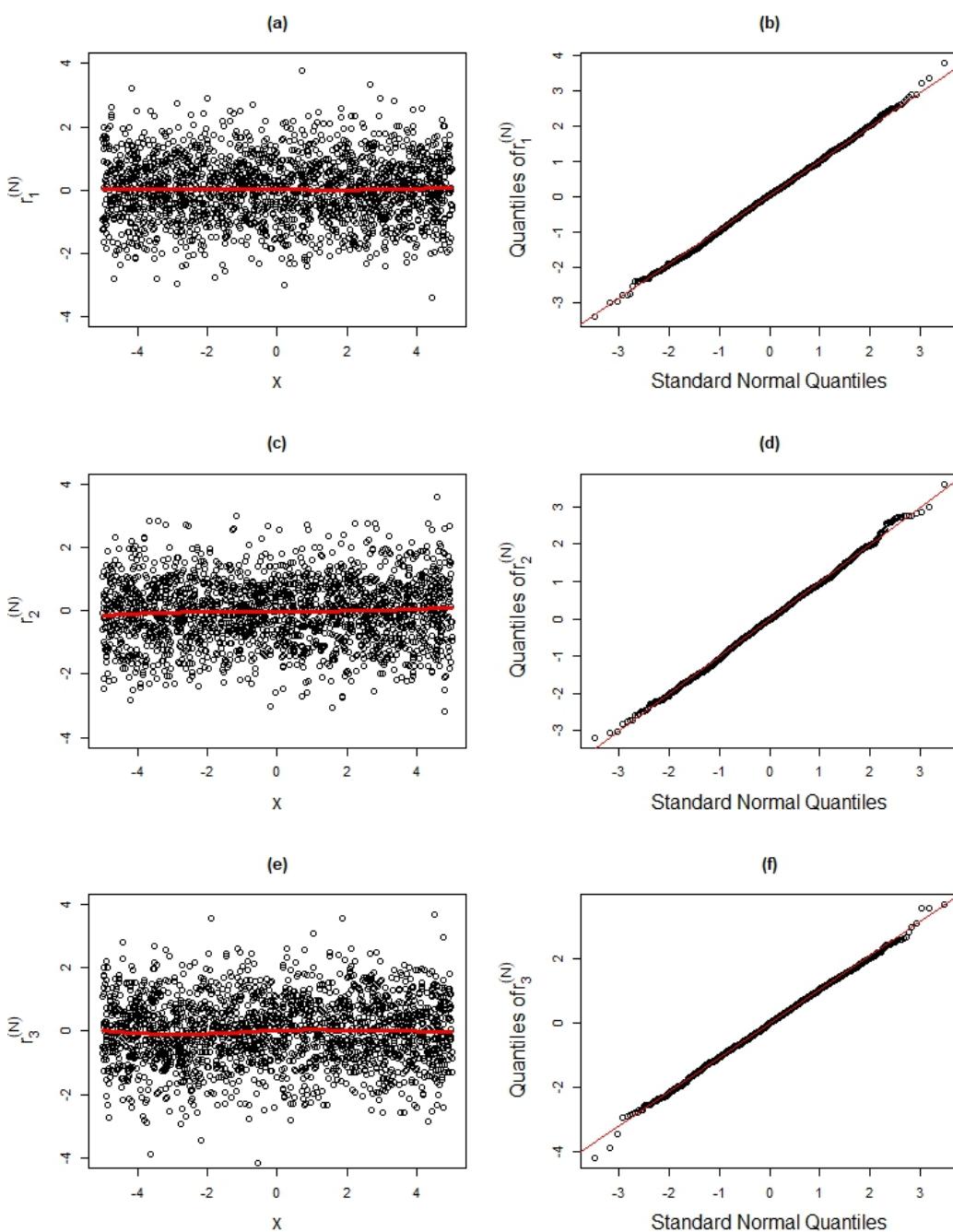


Figure S7: Model diagnostics for the improved model in Example 2 with the normalized surrogate residual defined in Section 3.4.1. Figures (a), (b) and (c) show the residuals-by-covariate plots of  $r_1^{(N)}$ ,  $r_2^{(N)}$  and  $r_3^{(N)}$  respectively; loess curves (red solid) are added as references. Figures (b), (d) and (f) are the corresponding QQ-plots (quantile-quantile plots) of above mentioned residuals.

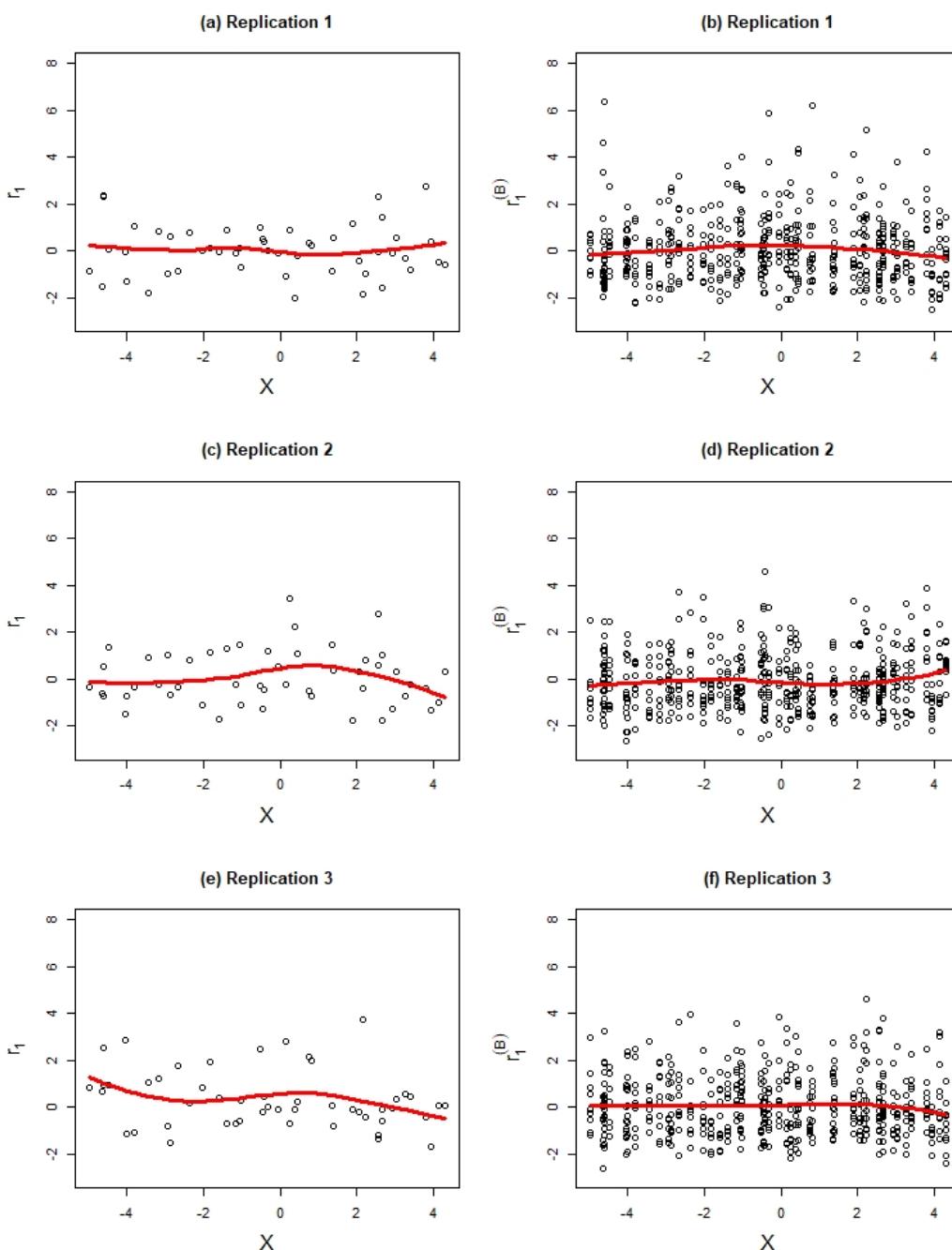


Figure S8: Three replications of  $r_1$ -by- $x$  plots by the surrogate residual ( $r_1$ ) and bootstrapping residual ( $r_1^{(B)}$ ) with 10 copies. Upper row: the first replication; Middle row: the second replication; Bottom row: the third replication.

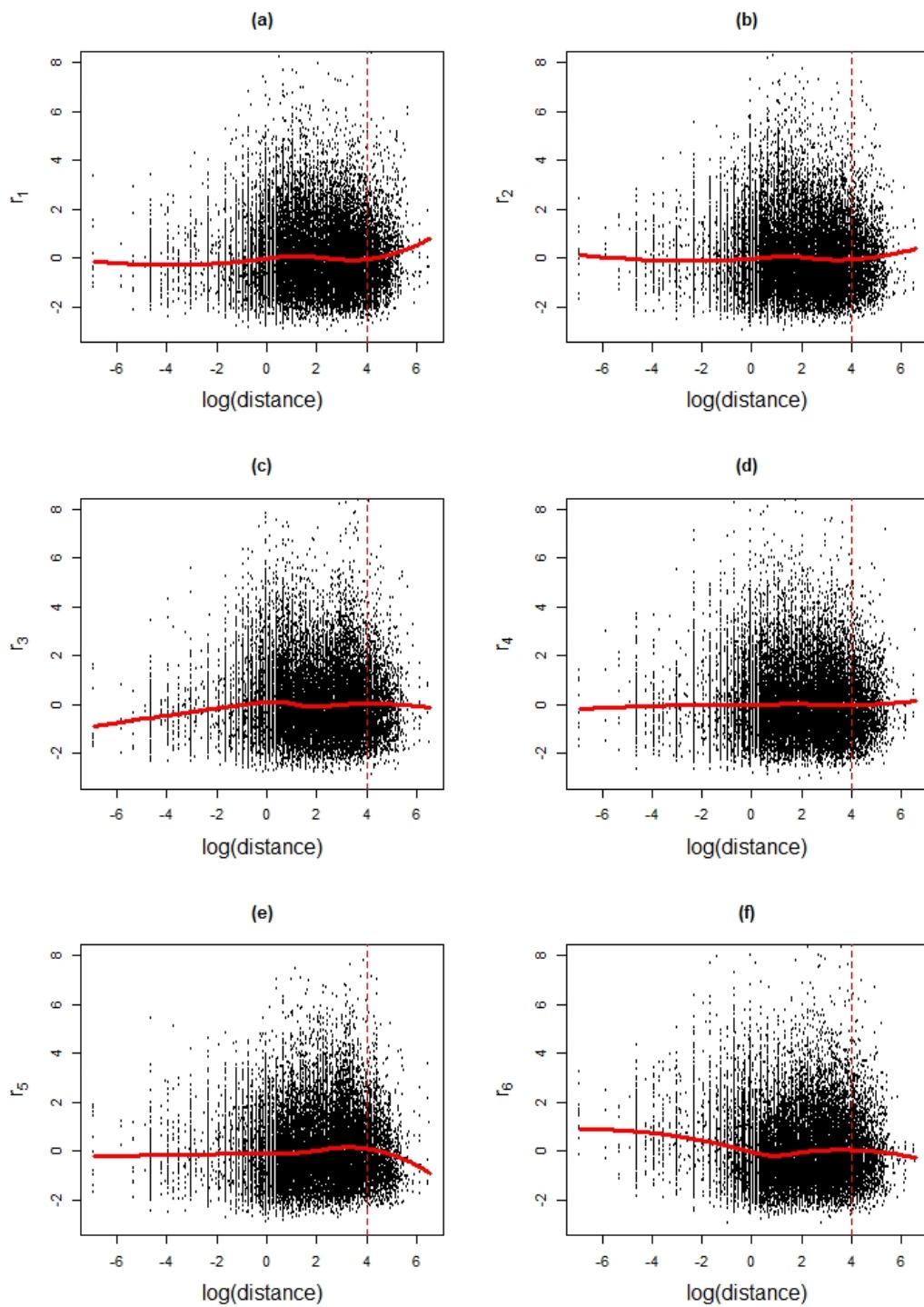


Figure S9: Residual-by- $\log(\text{distance})$  plots for our initial multinomial logistic regression fitted to the MPN data.  $r_1$  to  $r_6$  represent the surrogate residuals for the utilities of car as driver, car as passenger, bicycle, ebike, public transport, and walking, respectively. A vertical dashed line at  $X = 4$  and a loess curve (solid) are added for reference.

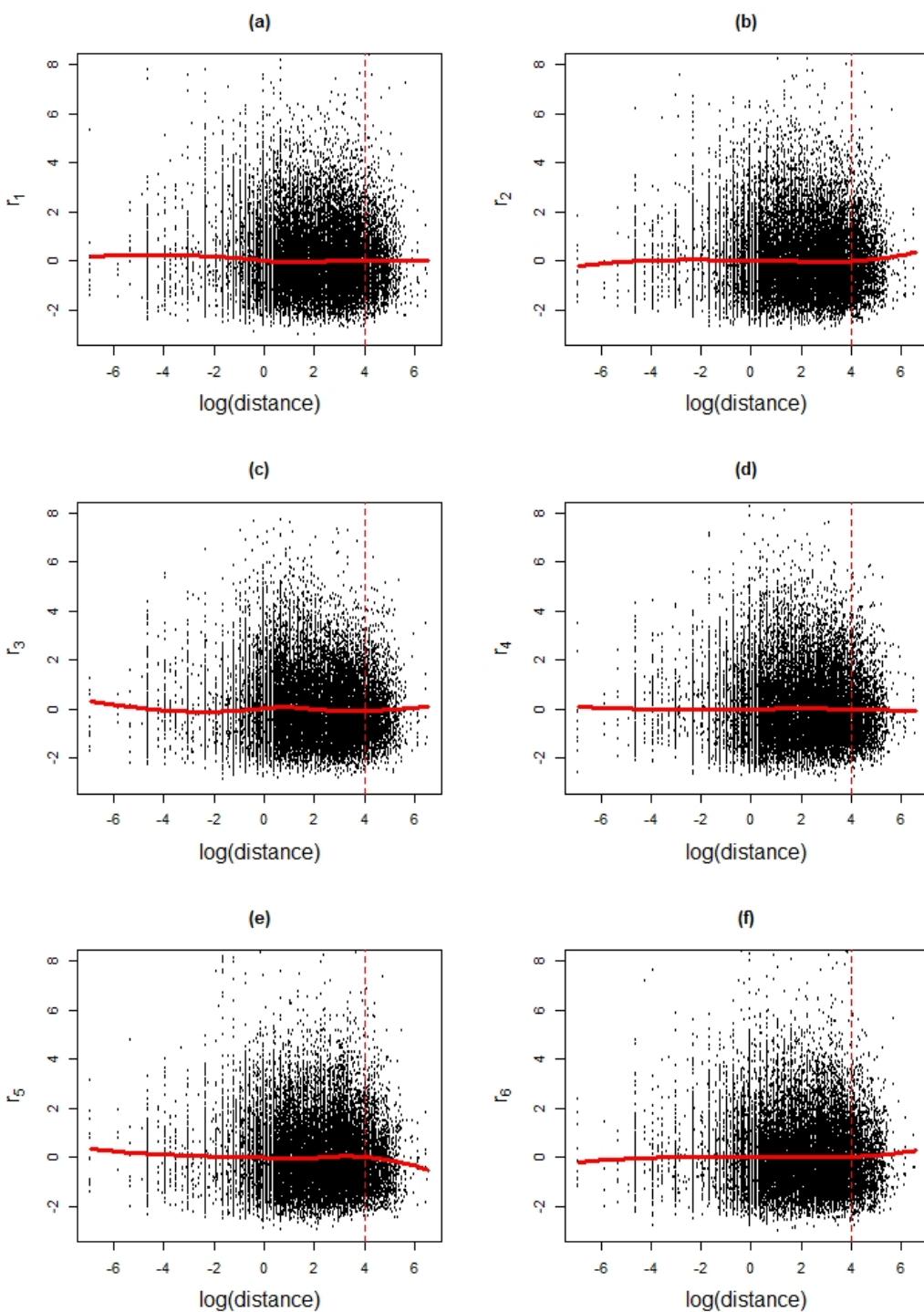


Figure S10: Residual-by- $\log(\text{distance})$  plots when natural logarithm term of distance are applied in the multinomial logistic regression fitted to the MPN data.  $r_1$  to  $r_6$  represent the surrogate residuals for the utilities of car as driver, car as passenger, bicycle, ebike, public transport, and walking, respectively. A vertical dashed line at  $X = 4$  and a loess curve (solid) are added for reference.

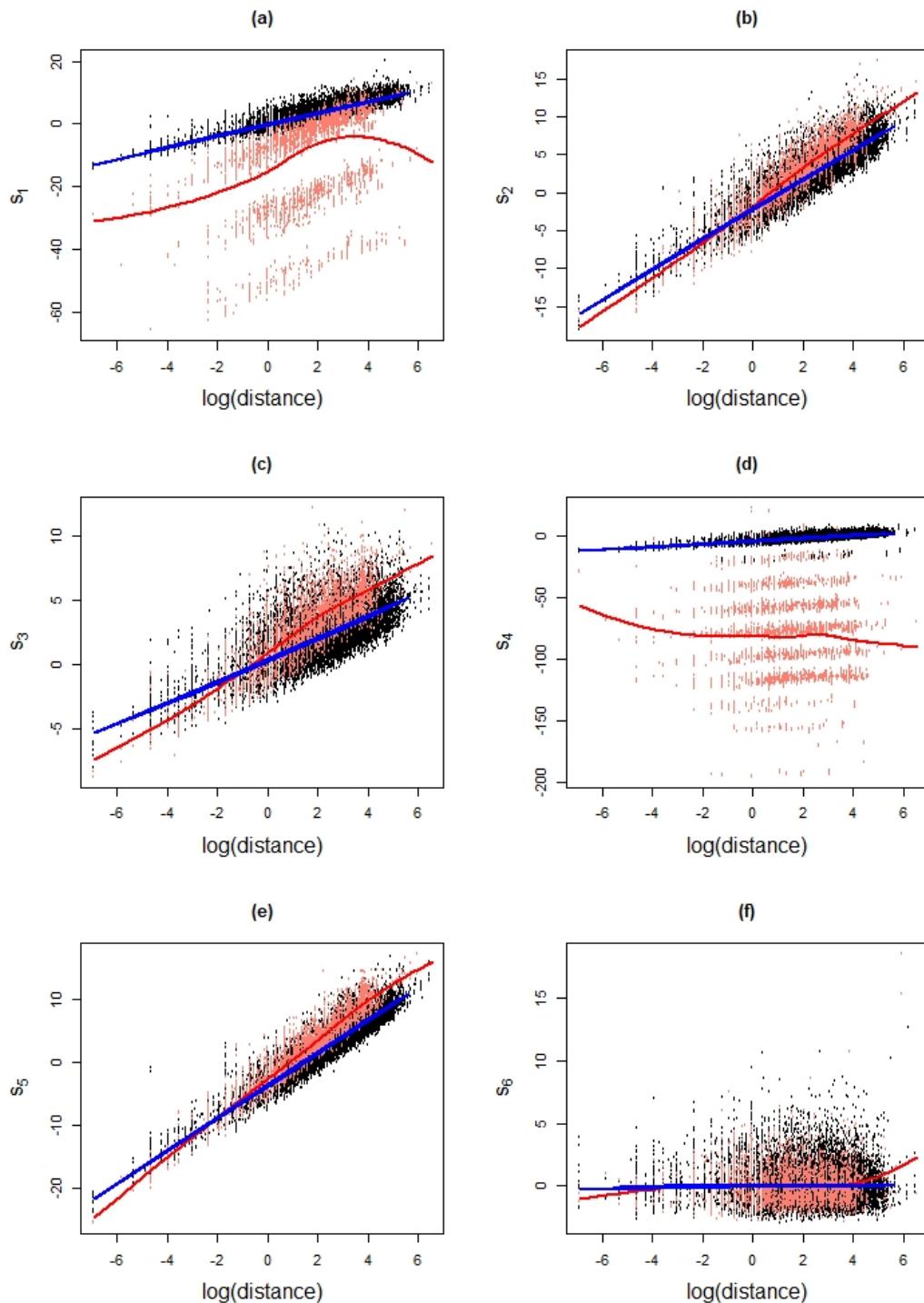


Figure S11: Plots of surrogate response versus  $\log(\text{distance})$  for the subjects who have a driving license and who do not have a driving license.  $s_1$  to  $s_6$  represent the surrogate responses corresponding to the utilities of car as driver, car as passenger, bicycle, ebike, public transport, and walking, respectively. The black points with a blue loess curve correspond to data from the subjects who have a driving license, while the orange points with a red loess curve correspond to data from the subjects who do not have a driving license.