

A peer-reviewed version of this preprint was published in PeerJ on 1 July 2014.

[View the peer-reviewed version](https://peerj.com/articles/467) (peerj.com/articles/467), which is the preferred citable publication unless you specifically need to cite this preprint.

Wang X. 2014. Modified generalized method of moments for a robust estimation of polytomous logistic model. PeerJ 2:e467
<https://doi.org/10.7717/peerj.467>

1 Modified generalized method of moments for a robust 2 estimation of polytomous logistic model

3 Xiaoshan Wang^{1,2*},

4 **1. Department of clinical and translational research/Forsyth Institute, Cambridge, MA,**
5 **USA**

6 **2. Department of Oral Health Policy and Epidemiology, Harvard School of Dental**
7 **Medicine, Cambridge, MA, USA**

8 * E-mail: xwang@forsyth.org

9 Abstract

10 The maximum likelihood estimation (MLE) method, typically used for polytomous logistic regression, is
11 prone to bias due to both misclassification in outcome and contamination in the design matrix. Hence,
12 robust estimators are needed. In this study, we propose such a method for nominal response data with
13 continuous covariates. A generalized method of weighted moments (GMWM) approach is developed for
14 dealing with contaminated polytomous response data. In this approach, distances are calculated based on
15 individual sample moments. And Huber weights are applied to those observations with large distances.
16 Mellow-type weights are also used to downplay leverage points. We describe theoretical properties of the
17 proposed approach. Simulations suggest that the GMWM performs very well in correcting contamination-
18 caused biases. An empirical application of the GMWM estimator on data from a survey demonstrates its
19 usefulness.

20 Introduction

21 Polytomous logistic regression models for multinomial data are a powerful technique for relating depen-
22 dent categorical responses to both categorical and continuous explanatory covariates [1, 2]. In practice,
23 however, the model building process can be highly influenced by peculiarities in the data. The maxi-
24 mum likelihood estimation (MLE) method, typically used for the polytomous logistic regression model
25 (PLRM), is prone to bias due to both misclassification in outcome and contamination in the design
26 matrix [3, 4]. Hence, robust estimators are needed.

27 For categorical covariates, we may apply MGP estimator [5], ϕ -divergence estimator [6], and robust
28 quadratic distance estimator [7]. The least quartile difference estimator can deal with overdispersion
29 problem [8]. But all these methods are difficult to be adapted for continuous covariates.

30 A generalized method of moments (GMM) estimation can be formed as a substitute of MLE. The
31 GMM is particularly useful when the moment conditions are relatively easy to obtain. GMM has been
32 extensively studied in econometrics [9–13]. Under some regularity conditions, the GMM estimator is
33 consistent [9]. With an appropriately chosen weight matrix, GMM achieves the same efficiency as the
34 MLE [14]. Furthermore, under certain circumstances, GMM provides more flexibility, such as dealing
35 with endogeneity through instrumental variables [15].

36 Like MLE, GMM estimation can be easily corrupted by aberrant observations [16]. Such observations
37 can bring up disastrous bias on standard parameter estimates if they are not properly accounted for,
38 see [17], [18], and [19]. So we propose a modified estimation method based on an outlier robust variant
39 of GMM. The method is different from the kernel-weighted GMM developed for linear time-series data
40 by [20] in that this is a data-driven method for defining weights. The new approach is evaluated using
41 asymptotic theory, simulations, and an empirical example.

42 The robust GMM estimator is motivated by the data from a 2006 study on hypertension in a sample
43 of the Chinese population. 520 people completed the survey. Observed variables included demographics,
44 social-economic status, weight, height, blood pressure, and food consumption. Sodium intakes were
45 calculated based on overall food consumption. Among those covariates, age, body mass index (BMI),
46 and sodium intakes are all continuous. Based on blood pressure measurements, subjects were classified
47 into 4 categories: Normal, Pre-hypertension, Stage 1 and Stage 2 hypertension. Table 1 lists the summary
48 statistics of the sample. One of the research objectives is to examine the association between hypertension

49 and risk factors in the population. Since the proportional odds assumption is violated (Score test for the
 50 proportional odds assumption gives $\chi^2 = 182.27$ with a degree of freedom of 8, $p < 0.0001$), we apply the
 51 polytomous logistic model, using the normal category as the reference level. In the case of J category,
 52 the polytomous logit model have $J - 1$ comparisons. Each comparison have a set of parameters for all
 53 covariates in the model. Therefore, the generalized logit model is not parsimonious when comparing with
 54 the proportional odds model. But the simultaneous estimation of all parameters is more efficient than
 55 separate models for each comparison. It is another option for ordinal response data, especially when a
 56 proportional odds model does not fit the data well. Table 2 lists the output from the model estimated
 57 by MLE. It is obvious that, if MLE is used, the estimates is inconsistent for sodium intakes, particularly
 58 the negative coefficient of sodium intake for the odds between the Stage 2 hypertension and the Normal
 59 categories. The inconsistency is more obvious when we plot the odds with respect to the sodium intake,
 60 the downward trend of the odds in Figure 2.A. This result contradicts the previous finding that there is
 61 a strong relationship between sodium intake and hypertension, see for example [21], [22] and references
 62 therein. Besides, Figure 2.A also shows another strange situation: the higher starting points for the odds
 63 between the Pre-hypertension and the Normal categories. The scatter plot (Figure 1) between distances
 64 and leverages suggests some observations are possible outliers: Observations 21, 33, 85, 92, 194, 274, 336,
 65 414, 459, 483, and 489 have large distances, which are blue-colored, and Observations 37, 83, 263, 459,
 66 483, 485, and 490 have large leverages, which are red-colored.

Table 1. Summary statistics for surveyed subjects.

Covariate		Hypertension categories			
		Normal	Pre-hypertension	Stage 1	Stage 2
Gender	Male	138	104	29	8
	Female	87	114	31	9
Age	Mean	43.2	48.8	54.3	60.3
	Std. Dev.	13.7	13.8	12.2	13.4
BMI	Mean	43.2	48.8	54.3	60.3
	Std. Dev.	13.7	13.8	12.2	13.4
Sodium Intake	Mean	3.7	3.7	4.6	2.7
	Std. Dev.	3.0	2.4	5.0	2.1

67 The paper is set up as follows. In the next section we presents the basic notations, model, and

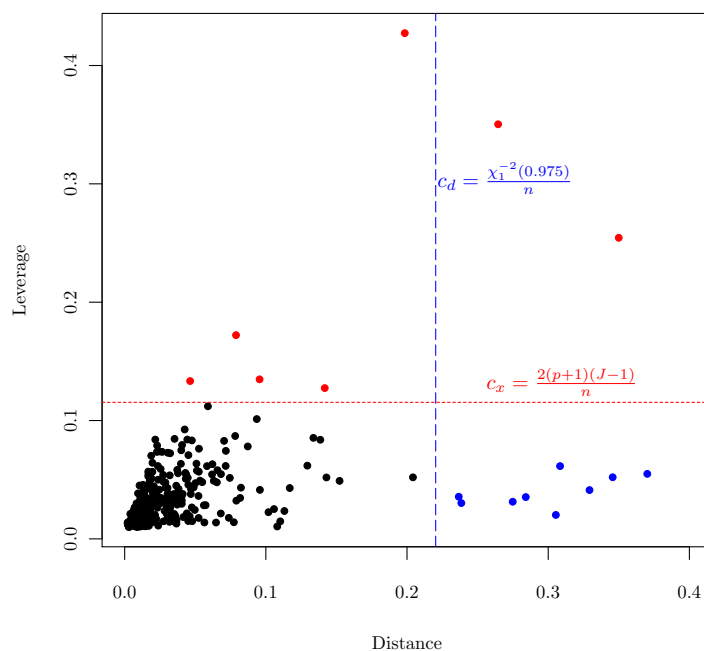


Figure 1. Scatter plot of distance vs. leverages, which are based on MLE. Criteria c_d for the distance and c_x for the leverage are demonstrated.

Table 2. Polytomous logistic regression of a hypertension data: coefficient estimates and standard errors from GMWM and MLE.

Var	Coefs	MLE			GMWM		
		Estimates	Std.Err	p value	Estimates	Std.Err	p value
Sex	β_{21}	0.7062	0.2022	0.0002	1.3339	0.2269	< 0.0001
	β_{31}	0.9789	0.3235	0.0012	1.0368	0.3013	0.0003
	β_{41}	1.4193	0.5746	0.0068	0.6753	0.2195	0.0010
Age	β_{22}	0.0350	0.0075	< 0.0001	0.0671	0.0086	< 0.0001
	β_{32}	0.0715	0.0121	< 0.0001	0.1139	0.0133	< 0.0001
	β_{42}	0.1096	0.0216	< 0.0001	0.0753	0.0103	< 0.0001
BMI	β_{23}	0.1147	0.0316	0.0001	0.1681	0.0360	< 0.0001
	β_{33}	0.2422	0.0474	< 0.0001	0.4382	0.0538	< 0.0001
	β_{43}	0.4351	0.0884	< 0.0001	0.2279	0.0388	< 0.0001
Sodium	β_{24}	0.0104	0.0349	0.3829	0.1831	0.0355	< 0.0001
	β_{34}	0.0919	0.0426	0.0155	0.2315	0.0486	< 0.0001
	β_{44}	-0.2699	0.1580	0.9562	0.2294	0.0353	< 0.0001

Std.Err = standard error.

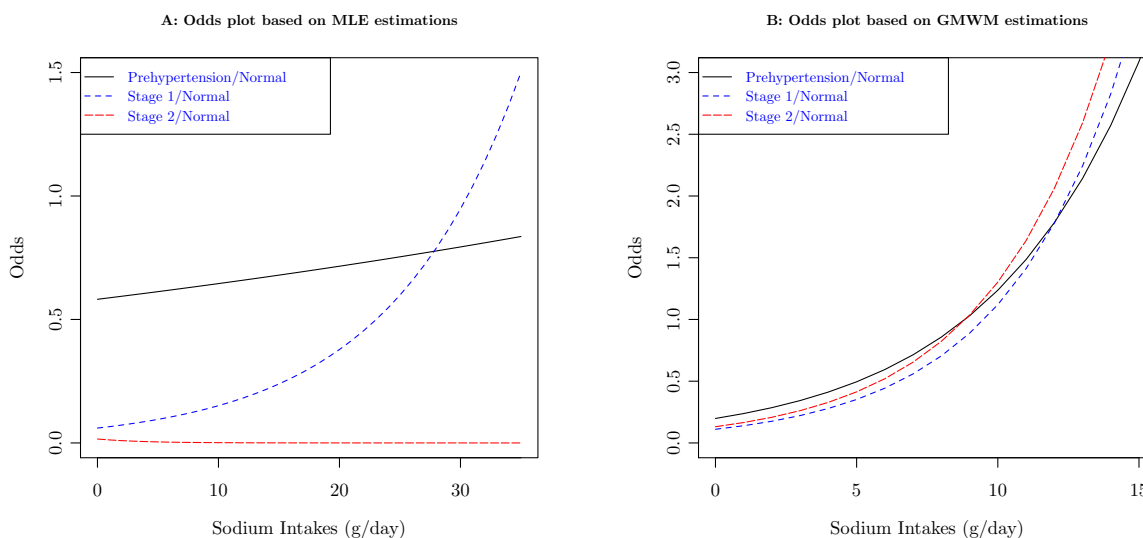


Figure 2. Compare odds plots of sodium intakes between MLE estimates and GMWM estimates on the population of Female, Age=40, and BMI=23.

68 standard GMM. Section “A robust GMM” introduces the outlier robust GMM estimator, and gives a
 69 detailed exposition of its implementation. In Section “Results”, we compares the performance of the
 70 standard MLE with the new estimator using a Monte-Carlo experiment. And we apply both estimators
 71 to real epidemiological data, and illustrate the usefulness of the robust estimator for application oriented
 72 researchers. We conclude with a discussion of advantages and limitations of the approach. The supporting
 73 document gathers the proofs of the asymptotic property.

74 Materials and Methods

75 The baseline-category logit model

76 Assume a random sample of size n from a large population. Each element in the population may be
 77 classified into one of J categories, denoted by $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{iJ})$ the multinomial trial for subject
 78 i , where $y_{ij} = 1$ when the response is in category j and $y_{ij} = 0$ otherwise, $i = 1, \dots, n$, $j = 1, \dots, J$.
 79 Thus, $\sum_j y_{ij} = 1$. Suppose p explanatory covariates, with at least one of them being continuous, are
 80 observed. Define $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ip})$, and $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$. We assume that $(\mathbf{y}_i, \mathbf{x}_i)$ are independently

81 and identically distributed (*i.i.d.*). Let $\pi_{ij} = \pi_j(\mathbf{x}_i) = P(Y_i = j|\mathbf{x}_i)$, denote the probability that the
 82 observation of Y belongs to category j , given covariates \mathbf{x}_i , we assume the relationship between the
 83 probability π_j and \mathbf{x} can be modeled as:

$$\log \left\{ \frac{\pi_j(\mathbf{x}_i)}{\pi_J(\mathbf{x}_i)} \right\} = \mathbf{x}_i^T \beta_j, \quad j = 2, \dots, J \quad (1)$$

84 where $\beta_j^T = (\beta_{j0}, \beta_{j1}, \dots, \beta_{jp})$. Here we set the first category as reference class. This model is called
 85 a baseline-category logit model [23] or generalized logit model [24]. MLE is usually used for obtaining
 86 parameter estimation of this model. Here we present an alternative estimation method formed with the
 87 GMM.

88 Estimation using GMM

89 The baseline-category logit model can be viewed as a multivariate model. Define $\mathbf{y}_i^{*T} = (y_{i2}, \dots, y_{iJ})$,
 90 since y_{i1} is redundant. Let $\mathbf{X}^T = (X_1^T, \dots, X_n^T)$ is a $n(J-1) \times (p+1)(J-1)$ matrix, with X_i^T , a
 91 $(J-1) \times (p+1)(J-1)$ matrix, defined as:

$$X_i^T = \begin{pmatrix} \mathbf{x}_i^T & & & \\ & \mathbf{x}_i^T & & \\ & & \dots & \\ & & & \mathbf{x}_i^T \end{pmatrix} \quad (2)$$

92 In the GMM framework, we define

$$u(\boldsymbol{\beta}) = X_i(\mathbf{y}_i^* - \boldsymbol{\pi}_i), \quad i = 1, \dots, n. \quad (3)$$

where $\boldsymbol{\pi}_i^T = (\pi_{i2}, \pi_{i3}, \dots, \pi_{iJ})$. And $\boldsymbol{\beta}^T = (\beta_2^T, \beta_3^T, \dots, \beta_J^T)$ is the $(p+1)(J-1)$ vector of unknown
 parameters. The population moment condition is

$$E\{u(\boldsymbol{\beta})\} = 0,$$

93 with the corresponding sample moment condition

$$U_n(\boldsymbol{\beta}) = \sum_{i=1}^n u(\boldsymbol{\beta}). \quad (4)$$

The GMM estimation of $\hat{\boldsymbol{\beta}}_M$ can be obtained by minimizing the following quadratic objective function

$$Q_n(\boldsymbol{\beta}) = U_n^T(\boldsymbol{\beta})\Sigma_n^{-1}(\boldsymbol{\beta})U_n(\boldsymbol{\beta}),$$

where $\Sigma_n(\boldsymbol{\beta})$ can be the empirical variance-covariance matrix given by

$$\Sigma_n(\boldsymbol{\beta}) = \frac{1}{n^2} \sum_{i=1}^n u^T(\boldsymbol{\beta})u(\boldsymbol{\beta}) - \frac{1}{n} U_n(\boldsymbol{\beta})U_n^T(\boldsymbol{\beta}).$$

94 Or, for the best efficiency of the GMM estimation, we can take the information matrix of the polytomous
95 logit model (PLRM), that is,

$$\Sigma_n(\boldsymbol{\beta}) = \sum_{i=1}^n X_i(D_i - \boldsymbol{\pi}_i\boldsymbol{\pi}_i^T)X_i^T. \quad (5)$$

96 where $D_i = \text{diagonal}(\boldsymbol{\pi}_i)$.

97 In general, $\hat{\boldsymbol{\beta}}_M$ can be computed via an iterative procedure [25]. Under standard regularity conditions,
98 the GMM estimator $\hat{\boldsymbol{\beta}}_M$ exists and converges in probability to the true parameter $\boldsymbol{\beta}_0$ [9]. A proof of
99 asymptotic normality of GMM can be found in Page 2148 of [13].

100 A robust GMM

101 In this section we introduce the outlier robust GMM estimator. In the following subsection, we specify
102 moment conditions used for robust estimation. And the details on the implementation of the estimator
103 follows.

104 The generalized method of weighted moments

105 The main principle used in the robust GMM estimator is that we replace moment conditions by a set of
106 observation weighted moment conditions. Instead of Equation (3), we define

$$u^w(\boldsymbol{\beta}) = w_i X_i(\mathbf{y}_i^* - \boldsymbol{\pi}_i) - c_i, \quad i = 1, \dots, n. \quad (6)$$

where $c_i = E\{w_i X_i(\mathbf{y}_i^* - \boldsymbol{\pi}_i)\}$. Then the estimation can be based on the moment conditions

$$E\{u^w(\boldsymbol{\beta})\} = 0.$$

107 Consequently, the generalized method of weighted moments (GMWM) estimates can be defined by

$$\hat{\boldsymbol{\beta}}^w = \arg \min_{\boldsymbol{\beta} \in \mathcal{B}} Q_n^w(\boldsymbol{\beta}). \quad (7)$$

108 where

$$Q_n^w(\boldsymbol{\beta}) = [U_n^w(\boldsymbol{\beta})]^T \{\Sigma_n^w(\boldsymbol{\beta})\}^{-1} U_n^w(\boldsymbol{\beta}), \quad (8)$$

109 with

$$U_n^w(\boldsymbol{\beta}) = \sum_{i=1}^n u^w(\boldsymbol{\beta}). \quad (9)$$

110 Here we take the summation as the sample moment condition. The advantage of using the summation is
111 that it can lead us to a direct estimation of covariance matrix.

112 It is clear to see that this definition is analogous to the standard GMM. If we choose $w_i = 1$ and
113 $c_i = 0$ for all observations, the moment conditions in (6) are reduced to the standard moment conditions.
114 Therefore, the standard GMM is a special case of the GMWM.

115 In order to specify the weights for the robust GMM estimator, we need the following definition of a
116 distance, which is based on individual moment conditions:

$$d_i(\boldsymbol{\beta}) = [u_i^w(\boldsymbol{\beta})]^T \{\Sigma_n^w(\boldsymbol{\beta})\}^{-1} u_i^w(\boldsymbol{\beta}), \quad i = 1, \dots, n. \quad (10)$$

117 The weight is assigned based on $d_i(\boldsymbol{\beta})$, that is, $w_d = w(d_i(\boldsymbol{\beta}))$. There are several alternative specifications
118 of weight functions available in the literature [17, 18]. In this study, the Huber's weights are applied:

$$w(d_i(\boldsymbol{\beta})) = \min\left(1, \frac{c_d}{d_i(\boldsymbol{\beta})}\right). \quad (11)$$

119 The above specification of weight function requires a value of the tuning constant c_d . Both the outlier
120 sensitivity and the efficiency of the estimator are determined by the constant. On the one hand, the
121 estimator should be reasonably efficient if the sample contains no outlier. On the other hand, the estimator

122 should be insensitive to outliers. To determine c_d , understanding the distribution of $d_i(\boldsymbol{\beta})$ is critical.
 123 Clearly, $u_i^w(\boldsymbol{\beta})$ is a column vector, and $d_i(\boldsymbol{\beta})$ is a scalar quadratic distance, so we set $c_d = \chi_1^{-2}(0.975)/n$,
 124 where $\chi_p^{-2}(\cdot)$ is the quantile of the χ^2 distribution with p degrees of freedom.

125 If we take the information matrix (5) of the PLRM as $\Sigma_n^w(\boldsymbol{\beta})$, we can computer a leverage for each
 126 observation:

$$H_i = X_i\{\Sigma_n^w(\boldsymbol{\beta})\}^{-1}X_i^T\sigma_i^w, \quad i = 1, \dots, n. \quad (12)$$

127 where σ_i^w is the i th component of $\Sigma_n^w(\boldsymbol{\beta})$. Then, a Mallows-type weight can be defined based on $trace(H_i)$;
 128 that is, $w_x = w(trace(H_i))$, to downplay the observations with high leverages. [26] suggest that the
 129 practical rule for isolating leverage points might set $c_x = 2(p+1)(J-1)/n$. In this study, we give
 130 observations with large leverages 0 weights,

$$w_x = w(trace(H_i)) = \begin{cases} 1 & \text{if } trace(H_i) \leq \frac{2(p+1)(J-1)}{n} \\ 0 & \text{otherwise.} \end{cases} \quad (13)$$

131 An approach often used to combine the two weights is $w_i = w_d \cdot w_x$. [27].

The consistency correction vector c_i is defined as

$$c_i = \left(w \left(d_i^{(1)}(\boldsymbol{\beta}) \right) - w \left(d_i^{(0)}(\boldsymbol{\beta}) \right) \right) / \text{diag} \left(\Sigma_n^w(\boldsymbol{\beta}) \right), \quad i = 1, \dots, n.$$

132 where $w \left(d_i^{(h)}(\boldsymbol{\beta}) \right) = w \left(X_i\{h - \pi_i(\boldsymbol{\beta})\} / \text{diag} \left[\Sigma_n^w(\boldsymbol{\beta}) \right]^{-1} \right)$ with $h = \{0, 1\}$, is the weight for y_i^* .

133 Implementation of the estimator

134 The continuous updating estimation method is applied in this study for estimating the regression coeffi-
 135 cients and corresponding variance. The procedure is detailed as follows:

- 136 1. Apply an initial value $\beta^{(0)}$ for computing $\Sigma_n(\boldsymbol{\beta})$.
- 137 2. Compute $d_i(\boldsymbol{\beta})$ using Equation (10) and H_i using Equation (12); assign weights correspondingly
 138 based on (11) and (13).
- 139 3. With the combined weights, calculate $\Sigma_n^w(\boldsymbol{\beta})$ and $U_n^w(\boldsymbol{\beta})$ in Equation (9).
- 140 4. Obtain the estimator $\hat{\boldsymbol{\beta}}_w^{(1)}$ by minimizing Q_n^w of Equation (8).

141 5. Go back to Step 1, replace $\beta^{(0)}$ with the estimator $\hat{\beta}_w^{(1)}$ in computing $\Sigma_n^w(\hat{\beta}_w^{(1)})$, and move to the
 142 next iteration.

143 6. Continue this procedure until convergence criteria are met.

144 For the starting value $\beta^{(0)}$, a reasonable choice is the MLE estimation based on the original data.

145 In the appendix, we proved that, under some regularity assumptions, we can have that $\hat{\beta}_w$ is consistent
 146 for β_0 . And by studying the behavior of the weighted moment equations in a neighborhood of β_0 ,
 147 we showed that the asymptotic linearity ensures the applicability of the central limit theorem for the
 148 asymptotic normality of GMWM.

149 Results

150 Monte Carlo simulations

In this section we investigate the properties of the GMWM estimator using a Monte-Carlo study. We generate data with three response categories and two covariates which are from multivariate normal distribution with 0 mean and identity covariance. The true coefficient matrix β_0 is

$$\beta_0 = \begin{pmatrix} \beta_{10} & \beta_{20} & \beta_{30} \\ \beta_{11} & \beta_{21} & \beta_{31} \\ \beta_{12} & \beta_{22} & \beta_{32} \end{pmatrix} = \begin{pmatrix} 0 & 1.0 & -0.3 \\ 0 & -0.8 & 0.7 \\ 0 & -1.0 & -0.5 \end{pmatrix}.$$

151 Based on the specified coefficient values and using the probability based on the model (1), we compute
 152 the category-specific probabilities for each subject. Then, using the computed probabilities, we deter-
 153 mine the most likely category to which each subject belongs. This decision is made through random
 154 generation from the multinomial distribution with the probability vector as a parameter. For instance,
 155 multinomial categories in R-Language are generated using *rmultinom*($n_i, N_i, \pi(\mathbf{x}_i)$) function, where
 156 $\pi(\mathbf{x}_i) = (\pi_1(\mathbf{x}_i), \dots, \pi_J(\mathbf{x}_i))$ is the probability vector, n_i is the number of random vectors to draw, and
 157 N_i is the total number of objects that are put into J -categories. In our case, $n_i = N_i = 1$ for all subjects
 158 and $J = 3$.

159 Two sample sizes, 100 and 1000, are examined. For each sample size, we run the simulation 1000
 160 times. Average biases and MSEs are calculated and tabulated. Table 3 shows the results from randomly

161 generated data with no outliers added. When the sample size is small, GMWM will give greater biases
 162 on β_{30} and β_{31} compared to the MLE method. For the sample size 1000, biases on these two parameters
 163 increase too, but not so obviously. Variances will also be inflated due to the weights we applied.

Table 3. Bias of parameter estimates and MSE from randomly generated data without outliers.

n	Parameter	True	MLE			GMWM		
			Bias	MSE	Coverage	Bias	MSE	Coverage
100	β_{20}	1.0	0.0666	0.1030	0.945	0.0488	0.1986	0.949
	β_{30}	-0.3	-0.0059	0.1206	0.957	-0.1440	0.5578	0.952
	β_{21}	-0.8	-0.0654	0.1190	0.938	-0.0513	0.2550	0.961
	β_{31}	0.7	0.0566	0.1892	0.963	0.2318	0.5468	0.923
	β_{22}	-1.0	-0.0853	0.1764	0.969	-0.0691	0.2380	0.950
	β_{32}	-0.5	-0.0624	0.1453	0.945	0.0203	0.3195	0.964
1000	β_{20}	1.0	0.0050	0.0087	0.956	0.0043	0.0181	0.962
	β_{30}	-0.3	-0.0055	0.0105	0.984	-0.0106	0.0333	0.950
	β_{21}	-0.8	-0.0039	0.0099	0.943	-0.0013	0.0251	0.956
	β_{31}	0.7	0.0081	0.0160	0.968	0.0162	0.0401	0.954
	β_{22}	-1.0	-0.0071	0.0145	0.987	-0.0025	0.0258	0.948
	β_{32}	-0.5	-0.0047	0.0122	0.948	0.0041	0.0361	0.947

164 Outliers are generated from a multivariate normal distribution with the mean vector = (2, 3) and
 165 identity covariance \mathbf{I}_2 . For these outliers, their responses are intentionally misclassified, that is, they are
 166 placed within a different category from those predicted categories based on the true parameters.

167 Table 4 lists simulation results with outliers added. For estimations from datasets with 5% outliers,
 168 bias correction from the GMWM is excellent. However, when the datasets have 10% outliers, biases
 169 on estimations of some parameters (β_{21} and β_{22} in this simulation) are decreased, but not completely
 170 corrected.

Table 4. Comparison between GMWM and MLE estimation from randomly generated data with outliers added.

Size	Parameter	5% contamination						10% contamination					
		GMWM			MLE			GMWM			MLE		
		Bias	MSE	Coverage	Bias	MSE	Coverage	Bias	MSE	Coverage	Bias	MSE	Coverage
100	β_{20}	0.0568	0.1102	0.956	0.0860	0.0884	0.957	0.0489	0.0999	0.971	0.0868	0.0819	0.970
	β_{30}	-0.0038	0.1427	0.954	-0.0055	0.1528	0.949	-0.0057	0.1510	0.945	-0.0431	0.1461	0.814
	β_{21}	-0.0392	0.1464	0.949	0.2377	0.1360	0.785	0.0319	0.1227	0.946	0.3607	0.1933	0.579
	β_{31}	0.0175	0.2020	0.944	-0.1072	0.1270	0.921	-0.0235	0.1770	0.943	-0.1631	0.1283	0.949
	β_{22}	0.0374	0.1207	0.949	0.3848	0.2115	0.578	0.0207	0.0968	0.945	0.6088	0.4151	0.526
	β_{32}	-0.0548	0.1572	0.956	-0.0964	0.0904	0.964	-0.0817	0.1349	0.977	-0.1069	0.0803	0.967
1000	β_{20}	0.0172	0.0189	0.939	0.0490	0.0102	0.932	0.0451	0.0202	0.944	0.0657	0.0120	0.900
	β_{30}	0.0012	0.0340	0.945	0.0124	0.0075	0.952	-0.0071	0.0336	0.952	-0.0111	0.0063	0.822
	β_{21}	0.0260	0.0242	0.937	0.2874	0.0885	0.101	0.0164	0.0207	0.936	0.3876	0.1545	0.002
	β_{31}	-0.0058	0.0356	0.950	-0.1423	0.0345	0.697	-0.0497	0.0346	0.917	-0.2269	0.0658	0.521
	β_{22}	0.0366	0.0237	0.936	0.4390	0.2032	0.000	0.0238	0.0182	0.938	0.6500	0.4322	0.000
	β_{32}	-0.0106	0.0292	0.951	-0.0538	0.0103	0.940	-0.0434	0.0250	0.953	-0.0629	0.0106	0.902

171 Application

172 For the hypertension data, the criterion for identifying observations with large distances is $c_d = 0.22$,
173 and the criterion for identifying leverage points is $c_x = 0.12$. Applying the GMWM estimator, those
174 blue-colored points in Figure 1 are automatically downweighted, and red-colored points have 0 weight.
175 The GMWM method indeed eliminates those inconsistencies: the coefficient of sodium intake for the
176 odds model between the Stage 2 hypertension and the Normal categories is no longer negative, see the
177 right side of Table 2.

178 As the results indicate, age, gender, and BMI all had significant impact on hypertension status. For
179 example, one unit increase in BMI resulted in an increase of 1.26 (95% confidence interval: 1.16 - 1.35)
180 times in likelihood to have Stage 2 hypertension when compared with the normal status. And with one
181 year age increase, a subject was 1.07 (95% CI: 1.06 - 1.10) times more likely to have Stage 2 hypertension
182 than to stay at the normal healthy status. Contrary to the MLE results for sodium intakes, which
183 were difficult to make a conclusion due to inconsistent estimate, we now find that sodium intakes were
184 statistically significant. When a daily intake of sodium increased one gram, a subject were 1.26 (95%
185 CI: 1.15 - 1.37) times more likely to have Stage 1 hypertension, and 1.25 (95% CI: 1.17 - 1.35) times
186 more likely to have Stage 2 hypertension. These results are consistent with the findings from previous
187 studies [21, 22].

188 Discussion

189 A reasonable choice to fit ordinal response data is the proportional odds model if the proportional odds
190 assumption is not violated. Proportional odds models can take the ordinal information into modeling.
191 And it reduces the number of parameters which is needed by the generalized logit model. Unfortunately,
192 our data does not met the fundamental assumption of proportional odds models, which makes us choose
193 to treat the outcome as a nominal response.

194 A datum with a nominal response and some continuous covariates is commonly seen in many scientific
195 areas, such as sociology, economy, and biomedical studies. In order to be able to deal with outliers,
196 we modified the GMM estimator to replace the standard moment conditions with weighted moment
197 conditions, so that aberrant observations automatically receive less weight. We proved that the proposed
198 method has good asymptotic behavior. When outliers are present, the GMWM estimator give much

199 smaller biases than the estimations derived from the traditional MLE method. This method can be
200 adapted to check whether results obtained with the traditional MLE approach are driven only by a few
201 outlying observations. The weights produced from the robust procedure can be used to diagnose the
202 cause of the differences and to indicate routes for model re-specification.

203 Appendix: Consistency and asymptotic normality

204 In this appendix, we introduce the assumptions for the asymptotic analysis of GMWM, and outline the
205 derivations on the main asymptotic properties of GMWM.

206 We make the following sets of regularity assumptions regarding properties of the moment functions
207 and identification assumptions.

208 *Assumption I*

209 I1. \mathcal{B} is a compact parametric space.

210 I2. Σ is a positive definite matrix.

I3. It holds that $E[u^w(\beta)] = 0$ if and only if $\beta = \beta_0$, and for any $\epsilon > 0$, that

$$\inf_{\beta \in \mathcal{B} \setminus \mathcal{N}(\beta_0, \epsilon)} \|E[u^w(\beta)]\| > 0$$

211 where $\mathcal{N}(\beta_0, \epsilon) = \{\beta \in \mathbb{R}^l \mid \|\beta - \beta_0\| < \epsilon\}$ is an open ϵ -neighborhood of a point β_0 .

212 *Assumption F*

213 F1. Let $u^w(\beta)$ be continuous in $\beta \in \mathcal{B}$, and be twice differentiable in β on $\mathcal{N}(\beta_0, \epsilon)$ almost surely.

214 F2. Expectation $E \sup_{\beta \in \mathcal{B}} \|u^w(\beta)\|$, $E \sup_{\beta \in \mathcal{N}(\beta_0, \epsilon)} \|\partial u^w(\beta)/\partial \beta_k\|$, and

215 $E \sup_{\beta \in \mathcal{N}(\beta_0, \epsilon)} \|\partial^2 u^w(\beta)/\partial \beta_k \partial \beta_l\|$ exists and are finite for $k, l = 1, \dots, p$.

216 *Assumption W*

217 W1. $\lim_{\epsilon \rightarrow 0} \sup_{\|\Delta\| \leq \epsilon} |w(\beta + \Delta) - w(\beta)| = 0$.

218 W2. $\lim_{\epsilon \rightarrow 0} \sup_{\|\Delta\| \leq \epsilon} |\partial w(\beta + \Delta)/\partial \beta - \partial w(\beta)/\partial \beta| = 0$.

219 When the above assumptions are met, we can prove that $\hat{\beta}_w$ is consistent for β_0 . We begin with studying
220 the behavior of the weighted moment equations in a neighborhood of β_0 . And proving their asymptotic
221 linearity is followed. The linearity ensures the applicability of the central limit theorem for the asymptotic
222 normality of GMWM.

Theorem 1. Let the Assumption F and I hold, then the GMWM estimator $\hat{\beta}_w$ is asymptotically normal, that is, $\sqrt{n}(\hat{\beta}_w - \beta_0) \xrightarrow{F} N(0, M^T S^w M)$ as $n \rightarrow \infty$, where

$$M = ((V^w)^T \Sigma V^w)^{-1} (V^w)^T \Sigma,$$

$$S^w = E \left[\frac{1}{n} \sum_{i=1}^n u^w(\hat{\beta}) u^w(\hat{\beta})^T \right].$$

223 with $V^w = E \left[\frac{\partial U^w(\hat{\beta})}{\partial \beta^T} \right]$.

224 We start with proving two lemmas before we present the proof of Theorem 1.

225 **Lemma 1.** Let the assumption F, I and W hold, and let $U_r^w(\beta)$ be the r th element of the vector $U^w(\beta)$,
226 $r = 1, \dots, p$. Then, for $0 < s < 1$,

$$\sup_{\|t\| \leq C} \left| \frac{1}{n} \sum_i \sum_{l=1}^p t_l \left\{ (\partial/\partial \beta_l) U_{i,r}^w \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial \beta_l) U_{i,r}^w(\beta) \right\} \right| = o_p(1). \quad (14)$$

227

Proof. For $l, r = 1, \dots, p$, by differentiating the i th component of $U_r^w(\beta)$, we get

$$\frac{\partial U_{i,r}^w(\beta)}{\partial \beta_l} = -w(\beta) X_i \frac{\partial \pi_i(\beta)}{\partial \beta_l} + \frac{\partial w_i(\beta)}{\partial \beta_l} X_i (y_i - \pi_i(\beta))$$

228 Then,

$$\begin{aligned} & \sup_{\|t\| \leq C} \left| \frac{1}{n} \sum_{i=1}^n \sum_{l=1}^p t_l \left\{ (\partial/\partial \beta_l) U_{i,r}^w \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial \beta_l) U_{i,r}^w(\beta) \right\} \right| \\ & \leq \frac{C}{n} \sum_{i=1}^n \sum_{l=1}^p t_l \sup_{\|t\| \leq C} \left| (\partial/\partial \beta_l) U_{i,r}^w \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial \beta_l) U_{i,r}^w(\beta) \right| \end{aligned}$$

229 and

$$\begin{aligned}
& \sup_{\|t\| \leq C} \left| (\partial/\partial\beta_l)U_{i,r}^w \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial\beta_l)U_{i,r}^w(\beta) \right| \\
& \leq \sup_{\|t\| \leq C} \left\{ \left| w_i \left(\beta + \frac{st}{\sqrt{n}} \right) - w_i(\beta) \right| \left| (\partial/\partial\beta_l)\pi_i \left(\beta + \frac{st}{\sqrt{n}} \right) \right| \right\} \\
& + \sup_{\|t\| \leq C} \left\{ \left| (\partial/\partial\beta_l)\pi_i \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial\beta_l)\pi_i(\beta) \right| |X_i w_i(\beta)| \right\} \\
& + \sup_{\|t\| \leq C} \left\{ \left| (\partial/\partial\beta_l)w_i \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial\beta_l)w_i(\beta) \right| \left| X_i \left(y_i - \pi_i \left(\beta + \frac{st}{\sqrt{n}} \right) \right) \right| \right\} \\
& + \sup_{\|t\| \leq C} \left\{ \left| \left(y_i - \pi_i \left(\beta + \frac{st}{\sqrt{n}} \right) \right) - \left(y_i - \pi_i(\beta) \right) \right| (\partial/\partial\beta_l)w_i(\beta) \right\}.
\end{aligned}$$

230 Then, by taking expectation at both sides,

$$\begin{aligned}
& E \left\{ \sup_{\|t\| \leq C} \left| (\partial/\partial\beta_l)U_{i,r}^w \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial\beta_l)U_{i,r}^w(\beta) \right| \right\} \\
& \leq \sup_{\|t\| \leq C} \left| w_i \left(\beta + \frac{st}{\sqrt{n}} \right) - w_i(\beta) \right| \sup_{\|t\| \leq C} \left| (\partial/\partial\beta_l)\pi_i \left(\beta + \frac{st}{\sqrt{n}} \right) \right| \\
& + \sup_{\|t\| \leq C} \left| (\partial/\partial\beta_l)\pi_i \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial\beta_l)\pi_i(\beta) \right| \sup_{\|t\| \leq C} |X_i w_i(\beta)| \\
& + \sup_{\|t\| \leq C} \left| (\partial/\partial\beta_l)w_i \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial\beta_l)w_i(\beta) \right| E \left\{ \sup_{\|t\| \leq C} \left| X_i \left(y_i - \pi_i \left(\beta + \frac{st}{\sqrt{n}} \right) \right) \right| \right\} \\
& + \sup_{\|t\| \leq C} \left| \left(y_i - \pi_i \left(\beta + \frac{st}{\sqrt{n}} \right) \right) - \left(y_i - \pi_i(\beta) \right) \right| \sup_{\|t\| \leq C} (\partial/\partial\beta_l)w_i(\beta).
\end{aligned}$$

Thus, by conditions F and W, we have

$$E \left\{ \sup_{\|t\| \leq C} \left| (\partial/\partial\beta_l)U_{i,r}^w \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial\beta_l)U_{i,r}^w(\beta) \right| \right\} \longrightarrow 0, \quad \forall i.$$

and

$$E \left\{ \sup_{\|t\| \leq C} \left| \frac{1}{n} \sum_{i=1}^n \sum_{l=1}^p t_l \left\{ (\partial/\partial\beta_l)U_{i,r}^w \left(\beta + \frac{st}{\sqrt{n}} \right) - (\partial/\partial\beta_l)U_{i,r}^w(\beta) \right\} \right| \right\} \longrightarrow 0, \quad \forall i.$$

231 Therefore, we have the results in (14). □

232 **Lemma 2.** *Let the Assumption F, I and W hold, it holds that*

$$\frac{1}{\sqrt{n}} \sup_{\|t\| \leq C} \left\| U_n^w(\beta_0 + n^{-\frac{1}{2}}t) - U_n^w(\beta_0) + V^w n^{-\frac{1}{2}}t \right\| = o_p(1), \quad (15)$$

233 as $n \rightarrow \infty$, where $V^w = E \left[\frac{\partial U^w(\beta)}{\partial \beta^T} \right]$.

Proof. Write

$$U_n^w(\beta_0 + n^{-\frac{1}{2}}t) - U_n^w(\beta_0) = \sum_{i=1}^n w_i(\beta_0 + n^{-\frac{1}{2}}t)u_i(\beta_0 + n^{-\frac{1}{2}}t) - \sum_{i=1}^n w_i(\beta_0)u_i(\beta_0)$$

234 By the Taylor expansion, $u_i(\beta_0 + n^{-\frac{1}{2}}t) = u_i(\beta_0) + n^{-\frac{1}{2}}t \left[\frac{\partial}{\partial \beta} u_i(\beta_0 + \frac{t}{\sqrt{n}}) \right]$, where $0 < s < 1$. Then, we
235 can write

$$\begin{aligned} & U_n^w(\beta_0 + n^{-\frac{1}{2}}t) - U_n^w(\beta_0) \\ &= \sum_{i=1}^n u_i(\beta_0) \left\{ w_i(\beta_0 + n^{-\frac{1}{2}}t) - w_i(\beta_0) \right\} \end{aligned} \quad (16)$$

$$+ \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i(\beta_0) t \frac{\partial}{\partial \beta} u_i(\beta_0) \quad (17)$$

$$+ \frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ w_i(\beta_0 + n^{-\frac{1}{2}}t) - w_i(\beta_0) \right\} t \frac{\partial}{\partial \beta} u_i(\beta_0) \quad (18)$$

$$+ \frac{1}{\sqrt{n}} \sum_{i=1}^n w_i \left(\beta_0 + \frac{t}{\sqrt{n}} \right) t \left\{ \frac{\partial u_i \left(\beta_0 + \frac{t}{\sqrt{n}} \right)}{\partial \beta} - \frac{\partial u_i(\beta_0)}{\partial \beta} \right\} \quad (19)$$

236 We will now show that terms (16), (18) and (19) are asymptotically negligible. As to the term (16), By
237 Assumption W1, $\{w_i(\beta_0 + n^{-\frac{1}{2}}t) - w_i(\beta_0)\} \rightarrow 0$, and $u_i(\beta_0)$ is independent of β . So we have the term
238 (16) tends to zero. Similarly, $\partial/\partial \beta u_i(\beta_0)$ is independent of β and t is bounded. Hence, the term (18)
239 tends to zero. Lemma 1 implies $\frac{\partial u_i(\beta_0 + \frac{t}{\sqrt{n}})}{\partial \beta} - \frac{\partial u_i(\beta_0)}{\partial \beta} \rightarrow 0$, as $n \rightarrow \infty$. So the term (19) can be neglect
240 too.

241 Now, let us analyze the term (17). Let $w_i^*(\beta_0)$ be the limit of $w_i(\beta_0)$. Rewrite (17) as

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n w_i(\beta_0) t \frac{\partial}{\partial \beta} u_i(\beta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n [w_i(\beta_0) - w_i^*(\beta_0)] t \frac{\partial}{\partial \beta} u_i(\beta_0) \quad (20)$$

$$+ \frac{1}{\sqrt{n}} \sum_{i=1}^n \left\{ w_i^*(\beta_0) t \frac{\partial}{\partial \beta} u_i(\beta_0) - E \left[w_i^*(\beta_0) t \frac{\partial}{\partial \beta} u_i(\beta_0) \right] \right\} \quad (21)$$

$$+ \frac{1}{\sqrt{n}} E \left[w_i^*(\beta_0) t \frac{\partial}{\partial \beta} u_i(\beta_0) \right] \quad (22)$$

The first term (20) is negligible because $\partial/\partial \beta u_i(\beta_0)$ is independent of β , t is bounded, and $[w_i(\beta_0) - w_i^*(\beta_0)] \rightarrow$

0. By the central limit theorem, each element of vector (21) converges in distribution to a normally distribution random variable with zero mean and a finite variance which is uniformly bounded by t . Hence, (21) is bounded in probability. The last term (20) is

$$\frac{1}{\sqrt{n}} E \left[w_i^*(\beta_0) t \frac{\partial}{\partial \beta} u_i(\beta_0) \right] = \frac{t}{\sqrt{n}} V^w$$

242 This proves the lemma. □

243 *Proof of Theorem 1.* Since $t_n = \sqrt{n} = \mathcal{O}_p(1)$ as $n \rightarrow \infty$ by Lemma 2, we can write (15) as

$$U_n^w(\beta_0 + n^{-\frac{1}{2}} t_n) - U_n^w(\beta_0) + V^w n^{-\frac{1}{2}} t_n = o_p(n^{-\frac{1}{2}}). \quad (23)$$

244 with a probability arbitrarily close to one uniformly in $t_n \in \{t : \|t\| \leq C\}$. Moreover, with $n^{-\frac{1}{2}} t_n = o_p(1)$,

245 $\partial U_n^w(\beta_0 + n^{-\frac{1}{2}} t_n) / \partial \beta \rightarrow V^w$ in probability as $n \rightarrow \infty$.

Note that the first order conditions of GMWM equal to 0, that is,

$$\frac{\partial Q_n^w(\beta_w)}{\partial \beta} = \left[\frac{\partial U_n^w(\beta_w)}{\partial \beta} \right]^T \Sigma(\beta_w) U_n^w(\beta_w) = 0,$$

246 Replace $U_n^w(\beta_w)$ with $U_n^w(\beta_0 + n^{-\frac{1}{2}} t_n)$ from Equation (23),

$$\begin{aligned} & \left[\frac{\partial U_n^w(\beta_0 + n^{-\frac{1}{2}} t_n)}{\partial \beta} \right]^T \Sigma(\beta_w) U_n^w(\beta_0 + n^{-\frac{1}{2}} t_n) \\ &= [V^w + o_p(1)]^T \Sigma(\beta_w) \left[U_n^w(\beta_0) - V^w n^{-\frac{1}{2}} t_n \right] = 0. \end{aligned}$$

247 Then we have

$$t_n = \sqrt{n}(\beta_w - \beta_0) = \sqrt{n} \left[(V^w)^T \Sigma(\beta_w) V^w \right]^{-1} (V^w)^T \Sigma(\beta_w) U_n^w(\beta_0) + o_p(1). \quad (24)$$

248 Next we examine the behavior of $\sqrt{n}U_n^w(\beta_0)$, which can be written as

$$\begin{aligned}\sqrt{n}U_n^w(\beta_0) &= n^{-\frac{1}{2}} \sum_{i=1}^n u_i(\beta_0)w_i(\beta_0) \\ &= n^{-\frac{1}{2}} \sum_{i=1}^n u_i(\beta_0) \{w_i(\beta_0) - w_i^*(\beta_0)\} \quad (25)\end{aligned}$$

$$+ n^{-\frac{1}{2}} \sum_{i=1}^n u_i(\beta_0)w_i^*(\beta_0). \quad (26)$$

249 Note that the term (25) is asymptotically negligible in probability due the the triangle inequality and
250 Assumption W1. The term (26) is stationary sequence of absolutely random variables. By Assumption
251 I3 and F2, (26) have zero mean and finite second moments. So the central limit theorem can be applied
252 on (26), giving $\sqrt{n}U_n^w(\beta_0) \sim N(0, S^w)$ [28, Section 25.3].

253 With Equation (24), we have asymptotic normality of β^w , and its asymptotic variance is given by
254 $M^T S^w M$ [28]. □

References

1. McCullagh P, Nelder JA (1989) Generalized Linear Models, Second Edition (Chapman & Hall/CRC Monographs on Statistics & Applied Probability). Chapman and Hall/CRC.
2. Liu I, Agresti A (2005) The analysis of ordered categorical data: An overview and a survey of recent developments. *TEST: An Official Journal of the Spanish Society of Statistics and Operations Research* 14: 1-73.
3. Pregibon D (1982) Resistant fits for some commonly used logistic models with medical applications. *Biometrics* 38: 485-498.
4. Copas JB (1988) Binary regression models for contaminated data. *Journal of the Royal Statistical Society Series B (Methodological)* 50: 225-265.
5. Victoria-Feser M, Ronchetti E (1997) Robust estimation for grouped data. *Journal of the American Statistical Association* 92: 333-340.
6. Gupta A, Kasturiratna D, Nguyen T, Pardo L (2006) A new family of bay estimators for polytomous logistic regression models based on ϕ -divergence measures. *Statistical Methods & Applications* 15: 159-176.
7. Flores, Garrido (2001) Robust logistic regression for insurance risk classification. Working Paper 01-64, Universidad Carlos III de Madrid.
8. Mebane J W R, Sekhon JS (2004) Robust estimation and outlier detection for overdispersed multinomial models of count data. *American Journal of Political Science* 48: 392-411.
9. Hansen LP (1982) Large sample properties of generalized method of moments estimators. *Econometrica* 50: 1029-1054.
10. Newey WK, West KD (1987) Hypothesis testing with efficient method of moments estimation. *International Economic Review* 28: 777-787.
11. Pakes A, Pollard D (1989) Simulation and the asymptotics of optimization estimators. *Econometrica* 57: 1027-1057.

12. Hansen LP, Heaton J, Yaron A (1996) Finite-sample properties of some alternative gmm estimators. *Journal of Business & Economic Statistics* 14: 262-280.
13. Newey WK, McFadden D (1986) Large sample estimation and hypothesis testing, In: *Handbook of Econometrics: vol. 4, 1st edn.* Elsevier.
14. Hayashi F (2000) *Econometrics.* Princeton University Press.
15. Baum CF, Schaffer ME, Stillman S (2002) Instrumental variables and gmm: Estimation and testing. *Boston College Working Papers in Economics* 545, Boston College Department of Economics.
16. Ronchetti E, Trojani F (2001) Robust inference with gmm estimators. *Journal of Econometrics* 101: 37-69.
17. Huber PJ (1981) *Robust Statistics (Wiley Series in Probability & Mathematical Statistics).* Wiley-Interscience.
18. Hampel FR, Ronchetti EM, Rousseeuw PJ, Stahel WA (1986) *Robust Statistics: The Approach Based on Influence Functions (Wiley Series in Probability and Statistics).* Wiley.
19. Rousseeuw PJ, Leroy AM (1987) *Robust Regression and Outlier Detection (Wiley Series in Probability and Statistics).* Wiley.
20. Kuersteiner GM (2012) Kernel-weighted GMM estimators for linear time series models. *Journal of Econometrics* 170: 399-421.
21. Institute of Medicine (2005) *Dietary Reference Intakes for Water, Potassium, Sodium, Chloride, and Sulfate.* The National Academies Press.
22. He FJ, MacGregor GA (2004) Effect of longer-term modest salt reduction on blood pressure. *Cochrane database of systematic reviews Online* 3: CD004937.
23. Agresti A (2002) *Categorical Data Analysis (Wiley Series in Probability and Statistics).* Wiley-Interscience.
24. Stokes ME, Davis CS, Koch GG (2001) *Categorical Data Analysis Using The SAS System.* WA (Wiley-SAS).

25. Hansen LP, Heaton J, Yaron A (1996) Finite-sample properties of some alternative gmm estimators. *Journal of Business & Economic Statistics* 14: 262-280.
26. Lesaffre E, Albert A (1989) Multiple-group logistic regression diagnostics. *Journal of the Royal Statistical Society Series C (Applied Statistics)* 38: 425-440.
27. Heritier S, Cantoni E, Copt S, Victoria-Feser MP (2009) *Robust Methods in Biostatistics* (Wiley Series in Probability and Statistics). Wiley.
28. Davidson J (1994) *Stochastic Limit Theory: An Introduction for Econometricians* (Advanced Texts in Econometrics). Oxford University Press, USA.