

# Weighted Robust Lasso and Adaptive Elastic Net Method for Regularization and Variable Selection in Robust Regression with Optimal Scaling Transformations

Tarek M. Omara

Department of Statistics, Mathematics and Insurance, Faculty of Commerce, Kafrelsheikh University, Kafrelsheikh, Egypt

**Abstract** In this paper, the weight least absolute deviation adaptive lasso optimal scaling method (WLAD-CATREG adaptive lasso) and weight least absolute deviation adaptive elastic net regression with optimal scaling method (WLAD-CATREG adoptive elastic net) will introduced, which is combined of weight least absolute deviation regression (WLAD-CATREG) and adaptive lasso (A-Lasso) or adaptive elastic net regression (A-Elastic net) with optimal scaling. Thus (WLAD-CATREG adoptive elastic net) method aim to automatically select variable, aspire to gropes effect and erase the bad effect of leverage points and outliers simultaneously, these aims cannot be achieved by (WLAD-CATREG), adaptive lasso regression (A-Lasso), weight robust adaptive lasso regression (WLAD-CATREG adoptive lasso), Weight least absolute deviation elastic net regression (WLAD-CATREG elastic net). Simulation study will be running to validated superiority of the (WLAD-CATREG adoptive Lasso) and (WLAD-CATREG adoptive elastic net).

**Keywords** A-Lasso, A-Elastic net, WLAD-CATREG, LAD-adoptive lasso, WLAD-CATREG adoptive lasso, WLAD-CATREG adoptive elastic net

## 1. Introduction

In regression models, if the variables are categorical variables, the relation between independent and dependent variables will be nonlinear. In this case, we used optimal scaling to transforming categorical variables to numeric variables, thus regression model becomes linear. This transform is done in simultaneous with the estimate parameters. The ordinary least squirt (OLS) is common method for estimating regression model but it sensitively to the outliers. In the case of regression model with optimal scaling transformations the effects of outliers can be as yet large (Peter (1993)), thus robust regression is suitable alternative. The robust regression include several methods one of them least absolute deviation (LAD) regression method which used to deal with outliers. This method developed to deal with leverage point (outlier in independent variables) by use weights which effect only on leverage point. In this side, weight least absolute deviation (WLAD) regression method proposed by (Giloni et al., 2006a, b, Olcay 2011). The regression model possibly is suffered from variable select problem, thus lasso regression method (Tibshirani (1996)) is appropriate because it does shrinkage

parameter and variable selection simultaneously. The same tuning parameter was used in lasso regression for all coefficient, so it suffered an palpable bias and not have the oracle properties. (Fan and Li (2001)), therefor adaptive lasso regression (A-Lasso) was used to allow several tuning parameters for several coefficients ((Zou (2006)), Hansheng et al. (2006)). To avoid outlier, adaptive lasso (A-Lasso) objective function been modified in to least absolute deviation adaptive lasso regression (LAD-adoptive lasso) which has oracle properties when we appropriately chosen tuning parameter (Hansheng et al. (2006) and Xu and Yin (2010)) and to avoid leverage points, (LAD-adoptive lasso) combined of weight least absolute deviation regression (WLAD) and adaptive lasso regression (A-lasso) (Olcay, 2011).

The lasso regression has limitation when  $p \gg n$ , and if there is a group of highly correlated covariates, the lasso select one variable form the group, therefore (Zou and Hastie (2005)) introduce elastic net regression which basically combined penalty of ridge and lasso ( $L_1$ ,  $L_2$ ) and has not an oracle properties. The elastic net does shrinkage parameter, variable selection and aspire to group effect simultaneously. To improve the elastic net (Samiran, 2007) combined of the adaptive lasso and elastic net and get adaptive elastic net which have the oracle property when  $p \gg n$ .

In this paper, for categorical regression model, we introduce new estimators (WLAD- CATREG adoptive elastic net) which considered appropriate way to deal with

\* Corresponding author:

tarek\_em@yahoo.com (Tarek M. Omara)

Published online at <http://journal.sapub.org/ajms>

Copyright © 2017 Scientific & Academic Publishing. All Rights Reserved

both leverage points and outliers, select variable and groping correlated variable simultaneously when  $p \gg n$ .

## 2. WLAD-CATREG Adoptive Lasso, WLAD-CATREG Adaptive Elastic Net

Consider the categorical regression model (CATREG)

$$v_r(y_i) = \sum_{j=1}^k v_j(x_{ij})\beta_j + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (1)$$

Where  $y$  is the discretized response variable,  $x_i$  is the discretized predictor variables,  $v_r$  is the function of transformation response variable,  $v_j$  is the function of transform predictor variables,  $\beta_j$  are the regression coefficients and  $\varepsilon$  is independent random error. The form of the last transformation based on optical scaling, so in the case of numerical scaling level the result of CATREG is the same of standardized multiple linear regression (Anita and van (2007)).

The (OLS) is the more commonly used for estimating the previous model. But this method is poorly when the model have outliers, thus many robust methods introduced. One important of the robust method called  $L_p$ CATREG regression methods. Since  $y_i$  and  $x_{ij}$  are standardized variables and if we fixed  $v_r(y_i)$ ,  $v_j(x_{ij})$  and  $\beta_j$  for all predictors  $l \neq j$  then  $L_p$  is written as

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n \left| v_r(y_i) - \sum_{j=1}^k v_j(x_{ij})\beta_j \right|^p \right], \quad p \geq 1 \quad (2)$$

In the special case, ( $p=1$ ), we get LAD CATREG or least absolute deviations regression with optimal scaling method

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n \left| v_r(y_i) - \sum_{j=1}^k v_j(x_{ij})\beta_j \right| \right] \quad (3)$$

The good leverage points, observation in the space of independent predictors, not affects significantly the LAD method, but bad leverage points, observation in the space of independent predictors but existing far from fit line, affects of it (Rousseeuw and Leroy (1987)). To deal with bad leverage points (Eills and Morgenthaler (1992)) suggested weight least absolute deviation (WLAD) regression method. An extinction of (WLAD) regression method, we get (WLAD-CATREG) regression method

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n w_i \left| v_r(y_i) - \sum_{j=1}^k v_j(x_{ij})\beta_j \right| \right] \quad (4)$$

Where  $0 < w_i \leq 1$ ,  $w_i$   $i=1, 2, \dots, n$  is the weights which will be chosen to breakdown leverage points. Chatterjee and Hadi, (1988) suggested weight depended of clean subset where

$$w_i = \sqrt{\frac{\min_j(h_j)}{h_i}}, \quad i=1, 2, \dots, n,$$

$$h_i = v(x^i)(v(X_C)^T v(X_C))^{-1} v(x^i)^T$$

$v(X_C)$  is clean subset,  $v(x^i)$  is the set of observation relative to clean subset. The discreet choice of weight lead to (WALD) estimator with fast computationally and high breakdown point. (Rousseeuw and Hubert (1997), (Olcay, 2011).) Used the robust distances (RD) to compute weights  $w_i$  which defined as

$$w_i = \min \left\{ 1, \frac{k}{[RD(x_i)]^2} \right\} = \min \left\{ 1, \frac{k}{\left[ \sqrt{(x_i - \hat{\mu})\hat{\Sigma}^{-1}(x_i - \hat{\mu})^T} \right]^2} \right\} \quad (5)$$

Where  $\hat{\mu}, \hat{\Sigma}$  are location and scatter estimators. Since robust distances (RD) identify leverage points increasing, (RD) lead to decrease weights  $w_i$  and thus leverage points corresponds to similar weights then it will be down weighted.

(Tibshirani (1996)) introduce lasso regression to combines estimate and variable selection. The lasso regression depend on minimize least squirt regression with the  $L_1$  norm condition. An extinction lasso regression, (Anita and van (2007)) introduce lasso penalties with (CATREG)

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n \left| v_r(y_i) - \sum_{j=1}^k v_j(x_{ij})\beta_j \right|^2 + \lambda_1 \sum_{j=1}^k |\beta_j| \right] \quad (6)$$

The lasso regression not have the oracle property, so ((Zou (2006)) introduce adaptive lasso (A-Lasso) which have oracle property when  $p \gg n$ . The adaptive lasso (A-Lasso) with (CATREG) defined as

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n \left| v_r(y_i) - \sum_{j=1}^k v_j(x_{ij})\beta_j \right|^2 + \lambda_1 \sum_{j=1}^k w_{s_j} |\beta_j| \right] \quad (7)$$

Where  $w_{s_j}, w_{s_l}$  is two weight vector,  $w_{s_j}, w_{s_l} > 0$ . In order to reach oracle property, ((Zou (2006), Olcay (2011)) define the weight vector as  $\hat{w}_{s_j} = |\hat{\beta}_j|^{-\gamma}$ ,  $j=1, 2, \dots, k$ , Where  $\gamma$  is a positive constant and  $\hat{\beta}_j$  is an elastic net estimator of  $\beta_j$ . (Zou and Zhang (2009)) used the other formula for weight vector as  $\hat{w}_{s_j} = |\hat{\beta}_j + 1/n|^{-\gamma}$ ,  $j=1, 2, \dots, k$  and (Wang et al (2007)) choose  $w_{s_j}$  as

$$\hat{w}_{s_j} = [\log(n)/n|\hat{\beta}_j|]^{-1} \quad j = 1, 2, \dots, k \quad (8)$$

which satisfies  $\sqrt{n} \hat{w}_{s_j} \rightarrow 0$  for  $j \leq k_0$  and  $\sqrt{n} \hat{w}_{s_j} \rightarrow \infty$  for  $j > k_0$  to avoid dividing zeros.

Since least squirt errors sensitive to outliers, (Hansheng and Chenlei. (2006)), introduced least absolute deviation

adaptive lasso regression (LAD-adoptive lasso). In this way, (Olcay (2011)) introduce weight least absolute deviation regression (WLAD) and adaptive lasso regression (A-lasso) to avoid leverage point add to the features of (LAD-adoptive lasso) which produce robust parameter with oracle property and select variables. The (WLAD-adoptive lasso) with (CATREG) (New) get as

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n w_i \left| v_r(y_i) - \sum_{j=1}^k v_j(x_{ij}) \beta_j \right| + n \sum_{j=1}^k \lambda_{1j} |\beta_j| \right] \quad (9)$$

On the other hand, (Zou and Hastie (2005)) introduce elastic net regression which add the group effect to the features of lasso regression. The elastic net regression depend on minimize least squirt regression with the  $L_1, L_2$  norm conditions. (Anita and van (2007)) introduce naive elastic net penalties with (CATREG) (CATREG naïve elastic net).

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n \left[ v_r(y_i) - \sum_{j=1}^k v_j(x_{ij}) \beta_j \right]^2 + \lambda_1 \sum_{j=1}^k |\beta_j| + \lambda_2 \sum_{j=1}^k \beta_j^2 \right] \quad (10)$$

(Samiran, 2007) introduce (Adaptive- elastic net) which have oracle property. We get (Adaptive- naive elastic net) with (CATREG) as

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n \left[ v_r(y_i) - \sum_{j=1}^k v_j(x_{ij}) \beta_j \right]^2 + \lambda_1 \sum_{j=1}^k w_{sj} |\beta_j| + \lambda_2 \sum_{j=1}^k \beta_j^2 \right] \quad (11)$$

Hong and Zhang (2010) developed the Adaptive naive elastic net to weight  $L_1$  and  $L_2$  penalty, thus we defined the adaptive elastic net with (CATREG) (New), in this case as

$$\beta = \operatorname{argmin} \left[ \frac{1}{2} \sum_{i=1}^n \left[ v_r(y_i) - \sum_{j=1}^k v_j(x_{ij}) \beta_j \right]^2 + n \lambda_1 \sum_{j=1}^k w_{sj} |\beta_j| + \frac{n}{2} \lambda_2 \sum_{j=1}^k w_{sj} \beta_j^2 \right] \quad (12)$$

The (LAD-CATREG adaptive elastic net (New)) which is scaled version of the (WLAD-CATREG naive adaptive elastic net (New)) defined as

$$\beta = \operatorname{argmin} \left[ \frac{1}{2} \sum_{i=1}^n \left[ v_r(y_i) - \sum_{j=1}^k v_j(x_{ij}) \beta_j \right]^2 + n \lambda_1 \sum_{j=1}^k w_{sj} |\beta_j| + \frac{n}{2} \lambda_2 \sum_{j=1}^k w_{sj} \beta_j^2 \right] \quad (13)$$

When  $p \gg n$ , to avoid leverage point and outliers in the elastic net regression, we defined (WLAD-CATREG naive adaptive elastic net (New)) as

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n w_i \left| v_r(y_i) - \sum_{j=1}^k v_j(x_{ij}) \beta_j \right| + n \sum_{j=1}^k \lambda_{1j} |w_{sj} \beta_j| + \lambda_2 \sum_{j=1}^k w_{sj} \beta_j^2 \right] \quad (14)$$

And we defined (WLAD-CATREG naive adaptive elastic net (New)) as

$$\beta_{\text{WLAD-CATREG naive adaptive elastic net}} = (1 + \lambda_2^{(n)}) \beta_{\text{WLAD-CATREG adaptive elastic net}} \quad (15)$$

### 3. Algorithm

To simplify (13), let

$$v_r(y_i)^* = w_i v_r(y_i), \quad v_r(\check{y}_i) = \begin{pmatrix} v_r(y_i)^* \\ 0 \end{pmatrix},$$

$$v_j(\check{x}_{ij}) = w_i v_j(x_{ij}), \quad \check{\beta}_j = \sqrt{(1 + \lambda_1)} \beta_j,$$

and define

$$(v_r(\check{y}_i), v_j(\check{x}_{ij})), \quad i=1,2,\dots,n,n+1,\dots,n+k, j=1,2,\dots,k.$$

where

$$(v_r(\check{y}_i), v_j(\check{x}_{ij})) = (v_r(y_i)^*, v_j(\check{x}_{ij}))$$

$$= (w_i v_r(y_i), w_i v_j(x_{ij})),$$

$$i=1,2,\dots,n, j=1,2,\dots,k,$$

$$(v_r(\check{y}_{n+j}), v_j(\check{x}_{(n+j)j})) = (0, n \lambda_{1j} e_j)$$

$e_j$  is a  $p$ -dimension vector with  $j$ th term equals one and all others equal to zero. We can rewrite (13) as a form of ridge as

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^{n+k} \check{W}_i \left| \begin{pmatrix} v_r(y_i)^* \\ 0 \end{pmatrix} - \sum_{j=1}^k \frac{1}{\sqrt{(1 + \lambda_1)}} \begin{pmatrix} v_j(\check{x}_{ij}) \\ w_{sj} \lambda_1 \end{pmatrix} \sqrt{(1 + \lambda_1)} \beta_j \right|^2 + \frac{\lambda_2}{(1 + \lambda_1)} \sum_{j=1}^k (1 + \lambda_1) |w_{sj} \beta_j|^2 \right]$$

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^{n+k} \check{W}_i \left| v_r(\check{y}_i) - \sum_{j=1}^k v_j(\check{x}_{ij}) \check{\beta}_j \right|^2 + \gamma_{ss} \sum_{j=1}^k |w_{sj} \check{\beta}_j|^2 \right]$$

Where  $\gamma_{ss} = \lambda_2 / (1 + \lambda_1)$ .

Define  $\check{X}_{ij_w} = \frac{\check{x}_{ij}}{w_{sj}}, \check{\beta}_{j_w} = w_s \check{\beta}_j$

Then

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^{n+k} \tilde{W}_i \left| v_r(\tilde{y}_i) - \sum_{j=1}^k v_j(\tilde{X}_{ij_w}) \tilde{\beta}_{j_w} \right| + \gamma_{ss} \sum_{j=1}^k |\tilde{\beta}_{j_w}|^2 \right] \quad (16)$$

$$\text{Where } \tilde{W}_i = \frac{1}{|v_r(\tilde{y}_i) - \sum_{j=1}^k v_j(\tilde{X}_{ij_w}) \tilde{\beta}_{j_w}|}$$

Before computing the WLAD-CATREG adaptive elastic net we must choose the weights  $w_i, w_{s_j}$  by use (and the tuning parameters  $\lambda_1, \lambda_2$  can be chosen Cross-validation (CV) on a two-dimensional but it would be computationally prohibitive (See Li and Jia (2010)).

So that we fixed  $\lambda_2$  and we used five-fold Absolute Cross-validation (ACV) to select tuning parameter  $\lambda_1$  which avoided leverage point. The Absolute Cross-validation (ACV) defined as

$$\text{ACV}(\lambda_2) = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{1 - H(\lambda_2)_{ii}} \right|$$

Where  $H_{ii}$  is hat matrix,  $H(\lambda_2) = X(X^T X + \lambda_2 I)^T X^T$ . We chosen weights  $w_i$  by the robust distances (RD) which defined as (5).

We computation the WLAD-CATREG adaptive elastic net by development the algorithm:

---

**Algorithm 1**

---

1. Input: Design Matrix  $X$ .
2. Find: The robust estimates  $\hat{\mu}, \hat{\Sigma}$ .
3. Calculate the initial regression estimator  $\hat{\beta}_j$  and the weight coefficient  $\hat{w}_s, \hat{w}$ .
4. Define  $r(y_i)^* = \hat{w}_{ir}(y_i)$ ,  $j(\tilde{x}_{ij}) = \hat{w}_{ij}(x_{ij})$ .
5. Solve the original WLAD-CATREG adaptive lasso

$$\beta = \operatorname{argmin} \left[ \sum_{i=1}^n w_i \left| r(y_i) - \sum_{j=1}^k j(x_{ij}) \beta_j \right| + n \sum_{j=1}^k \lambda_{1j} |\beta_j| \right]$$

and the WLAD-CATREG naive adaptive elastic net

$$\hat{\beta} = \operatorname{argmin} \left[ \sum_{i=1}^{n+k} \hat{w}_i \left| r(\tilde{y}_i) - \sum_{j=1}^k j(\tilde{X}_{ij_w}) \tilde{\beta}_{j_w} \right| + \gamma_{ss} \sum_{j=1}^k |\tilde{\beta}_{j_w}|^2 \right]$$

6. Calculate WLAD-CATREG adaptive elastic net

$$\hat{\beta}_{\text{WLAD-CATREG adaptive elastic net}} = (\mathbf{1} + \lambda_2^{(m)}) \hat{\beta}$$

Output:

$\hat{\beta}_{\text{WLAD-CATREG adaptive lasso}}$

$\hat{\beta}_{\text{WLAD-CATREG naive adaptive elastic net}}$ ,

$\hat{\beta}_{\text{WLAD-CATREG adaptive elastic net}}$

for  $j=1, 2, \dots, k$

---

## 4. The Properties of Estimators

In this section, we discuss the asymptotic properties of the WLAD-CATREG adaptive adaptive elastic net. At the first, we rooting some convenience and definitions. We decompose the  $\beta$  vector as

$$\beta = (\beta_a^T, \beta_b^T)^T = [(\beta_1, \beta_2, \dots, \beta_{k_0})^T, (\beta_{k_0+1}, \dots, \beta_k)^T]^T,$$

decompose the predictor variables  $x_i$  as

$$(v_j(x_{ia}), v_j(x_{ib})) = [(v_j(x_{i1}), v_j(x_{i2}), \dots, v_j(x_{ik_0}))^T, (v_j(x_{i(k_0+1)}), \dots, v_j(x_{ik}))^T]^T$$

Defined  $a_n = \max[\lambda_j : 1 \leq j \leq k_0]$  and  $b_n = \min[\lambda_j : k_0 < j \leq k]$  where  $\lambda_j$  is a function of  $n$ . (Wang and Leng (2007)).

Let  $\hat{\beta}_w = (\hat{\beta}_{ewa}^T, \hat{\beta}_{ewb}^T)^T$  be identical WLAD-CATREG adoptive net elastic estimator. Consider the linear regression model (1) with independent and assume the following conditions:

A1: The identically error with median zero and cumulative distribution function  $F$  which is positive and continues,

A2: The covariance of predictor variables  $x_i$  exists and positive definite

A3:  $W$  is  $n \times n$  diagonal matrix with known positive value ( $w_i, i=1, 2, \dots, n$ ),  $\max w_i = O(1)$  and  $\max w_i^{-1} = O(1)$

$$\text{A4: } \lim_{n \rightarrow \infty} \frac{X^T W X}{n} = \left( Q_n + \frac{\gamma_s^{(n)}}{n} I \right) \rightarrow Q$$

where  $Q$  is a positive definite.

$$\text{A5: } \sqrt{n} \lambda_1 = o(1), \sqrt{n} \lambda_2 = o(1)$$

$$\text{A6: } \lim_{n \rightarrow \infty} \frac{\gamma_{ss}}{\sqrt{n}} \sqrt{\sum_{j \in \mathcal{A}} |\tilde{\beta}_{j_w}|^2} = 0$$

The assumes (A1, A3, A4) are the same assumes for Olcay (2011), the assume (A4) is the (A6) assume for (Zou and Zhang (2009)) and the assumes A1, A2 the same assume for (Pollard, 1991).

**Lemma (1):** For the model (1), if it satisfies Assumptions A1: A6, then LAD-CATREG adaptive elastic net

$\hat{\beta}^T = (\hat{\beta}_a^T, \hat{\beta}_b^T)^T$  (14) must satisfies the following:

$$P(\hat{\beta}_b = 0) \rightarrow 1$$

$$\sqrt{n}(\hat{\beta}_a^T - \beta_a) \rightarrow N(0, \Sigma_{w_0}^{-1}/4f^2(0))$$

Proof lemma (1):

$$\text{Let } \tilde{\beta}_{j_w} = \hat{\beta}_{j_w} + \frac{\tau_j}{\sqrt{n}}$$

Then, we can rewrite (14) as

$$\Psi(\tau) = \sum_{i=1}^{n+k} \left| \tilde{y}_i - \sum_{j=1}^k \tilde{X}_{ij_w} \left( \hat{\beta}_{j_w} + \frac{\tau_j}{\sqrt{n}} \right) \right| + \gamma_{ss} \sum_{j=1}^k \left| \left( \hat{\beta}_{j_w} + \frac{\tau_j}{\sqrt{n}} \right) \right|^2$$

$$\tau_j = \sqrt{n}(\tilde{\beta}_{j_w} - \hat{\beta}_{j_w})$$

$$\Psi(0) = \sum_{i=1}^{n+k} \left| \dot{y}_i - \sum_{j=1}^k \ddot{X}_{ij_w} \tilde{\beta}_{j_w} \right| + \gamma_{ss} \sum_{j=1}^k |\tilde{\beta}_{j_w}|^2$$

Let  $D_n(\tau) = \Psi(\mu) - \Psi(0) =$

$$\begin{aligned} & \sum_{i=1}^{n+k} \left| \dot{y}_i - \sum_{j=1}^k \ddot{X}_{ij_w} \left( \tilde{\beta}_{j_w} + \frac{\tau_j}{\sqrt{n}} \right) \right| - \sum_{i=1}^{n+k} \left| \dot{y}_i - \sum_{j=1}^k \ddot{X}_{ij_w} \tilde{\beta}_{j_w} \right| \\ & + \gamma_{ss} \sum_{j=1}^k \left| \left( \tilde{\beta}_{j_w} + \frac{\tau_j}{\sqrt{n}} \right) \right|^2 - \gamma_{ss} \sum_{j=1}^k |\tilde{\beta}_{j_w}|^2 \\ & = \sum_{i=1}^{n+k} \left\{ \left| \dot{y}_i - \sum_{j=1}^k \ddot{X}_{ij_w} \left( \tilde{\beta}_{j_w} + \frac{\tau_j}{\sqrt{n}} \right) \right| - \left| \dot{y}_i - \sum_{j=1}^k \ddot{X}_{ij_w} \tilde{\beta}_{j_w} \right| \right\} \\ & + \gamma_{ss} \sum_{j=1}^k \left\{ \left| \left( \tilde{\beta}_{j_w} + \frac{\tau_j}{\sqrt{n}} \right) \right|^2 - |\tilde{\beta}_{j_w}|^2 \right\} \\ & \geq \sum_{i=1}^{n+k_0} \left\{ \left| \dot{y}_i - \sum_{j=1}^{k_0} \ddot{X}_{ij_w} \left( \tilde{\beta}_{j_w} + \frac{\tau_j}{\sqrt{n}} \right) \right| - \left| \dot{y}_i - \sum_{j=1}^{k_0} \ddot{X}_{ij_w} \tilde{\beta}_{j_w} \right| \right\} \\ & + n a_n \sum_{j=1}^{k_0} \{ |\tau_j|^2 \} \\ & \geq Z_n^{(1)}(\tau) + Z_n^{(2)}(\tau) \end{aligned} \tag{A1}$$

Knight (1998) holds that for  $x \neq 0$

$$\begin{aligned} |x - y| - |x| &= -y[I(x > 0) - I(x < 0)] \\ &+ 2 \int_0^y [I(x \leq s) - I(x \leq 0)] ds \end{aligned}$$

Using this equation, the first item at (A1)  $Z_n^{(1)}(\tau)$  can be expressed as

$$\begin{aligned} & -n^{-\frac{1}{2}} \sum_{i=1}^n \sum_{j=1}^k \tau_j \ddot{X}_{ij_w} [I(\dot{\epsilon}_i > 0) - I(\dot{\epsilon}_i < 0)] \\ & + 2 \sum_{i=1}^n \sum_{j=1}^k \int_0^{n^{-\frac{1}{2}} \mu_j \ddot{X}_{ij_w}} [I(x \leq s) \\ & - I(x \leq 0)] ds \end{aligned} \tag{A2}$$

Using the central limit theorem, the first item converges in distribution to  $\mu^T H$ , where  $\mu^T = (\mu_1, \mu_2, \dots, \mu_k)^T$  and  $H$   $p$ -dimensional normal random vector with mean 0 and variance matrix

$$\begin{aligned} \Sigma_w &= \text{cov}(\ddot{X}_w). \text{ Since } \lambda_1 = o(n^{-1/2}), \\ \max w_i &= O(1) \text{ and } \max w_i^{-1} = O(1) \text{ then } \Sigma_w = X^T W^2 X. \end{aligned}$$

Follow from the proofs of lemma1 (Wang et al. (2007)), the second item converges to

$$\begin{aligned} f(0)\tau^T \Sigma_w \tau &= f(0)\tau^T X^T W^2 X \tau \\ Z_n^{(1)}(\tau) &\rightarrow \tau^T H + f(0)\tau^T X^T W^2 X \tau \end{aligned}$$

Since  $\sqrt{n}a_n \rightarrow 0$ , the second item at (A1)  $Z_n^{(2)}(\tau)$  converges to 0 in probability.

Then

$$D_n(\tau) \rightarrow \mu^T H_0 + f(0)\tau^T X^T W^2 X \tau$$

Where  $H_0$   $p$ -dimensional normal random vector with mean 0 and variance matrix  $\Sigma_{w_0}$ .

Follow from the proofs of lemma1 (Wang et al. (2007)).

### 5. The Simulation Study

In this section, the performance of the WLAD-CATREG adaptive elastic net estimates are examined via the simulation. We simulate data sets from the true model  $v_r(y_i) = v(x_i)^T \beta_j + \sigma \epsilon_i$  where  $v(x_i) = [v(x_{i1}), v(x_{i2}), \dots, v(x_{ik})]^T$ . To ensure the situation of grouping variable, the simulated date consist of a training set, an independent validation set and an independent test set. In this simulate, we use the 50 simulated data sets each consisting of (30 training set/30 independent validation set /100 independent test set) observations and have 30 variables, 10 categorical variables each of them containing nine category and 20 numerical variables. The numerical variables generate  $X_{2i} \sim N_k(2, I), i = 1, 2, \dots, m$  recalling to the model  $v_r(y_{2i}) = \sum_{j=1}^k v_j(x_{2ij}) \beta_{2j}$  where  $\beta_{1j} \neq \beta_{2j}$  and the categorical variables generate by Bernoulli distribution with nine category. The correlation between pairs of variables  $r = \text{corr}(X_i, X_j), i \neq j$  was taken (0.20, 0.60). We have used WLAD-CATREG-1 algorithm to compute the LAD-adaptive lasso estimator and WLAD-CATREG adaptive lasso estimator and WLAD-CATREG-EN algorithm to compute the LAD-CATREG adaptive elastic net estimator and WLAD-CATREG adaptive elastic net estimator. To contaminate the data, we generated the contamination rate ( $\tilde{\epsilon} = 20\%, 40\%$ ). We will choose the place of outliers from observations on are works well. And the elastic net-CATREG methods select more variables than the lasso methods. In adding, the WLAD naïve adoptive elastic net-CATREG make improvements to a number of selected variables. categorical variable randomly and replacing the category on the variables by extreme values (such as 1 or 9) and will generate the outliers on numerical variable by  $v(x_{1i}) \sim N_k(0, 0.5^{|i-j|}), i = 1, 2, \dots, n - m$  where  $i \neq j$  recalling to the model  $v_r(y_{1i}) = \sum_{j=1}^k v_j(x_{1ij}) \beta_{1j} + \sigma \epsilon_i$ .

To ensure the achievement the normality of error with the possibility of the existence of outlier of  $y$ , we generate  $\epsilon$  from  $t$  distribution with two degrees of freedoms  $\epsilon \sim t(2)$  and set  $\sigma = 3$ . We let

$$\beta_1 = (1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 2, 2, 2, 2, 2, 3, 3, \underbrace{3, 3, 3, 3, 1, 1, 1, 1, 1, 1, 1, 1}_{})$$

and for the weight vector, we assume  $\gamma = 5$  Since elastic parameters have two tuning parameters, we used 2-dimensional cross-validation. For this method, we use  $\lambda_2 = 0.2$ .

**Table (1).** ACV( $\lambda_2$ ) and number of non-zero coefficients for the simulation study

*Estimators	ACV( $\lambda_2 = 0.2$ )				Number of non-zero coefficients			
	$r = 0.2$		$r = 0.6$		$\check{\epsilon}=20\%$		$\check{\epsilon}=40\%$	
	$\check{\epsilon}=20\%$	$\check{\epsilon}=40\%$	$\check{\epsilon}=20\%$	$\check{\epsilon}=40\%$	$\check{\epsilon}=20\%$	$\check{\epsilon}=40\%$	$\check{\epsilon}=20\%$	$\check{\epsilon}=40\%$
(1)	62.6	60.3	63.6	62.3	---	---	---	---
(2)	61.3	60.2	61.5	59.8	---	---	---	---
(3)	63.5	61.4	63.5	60.8	12	13	12	11
(4)	60.6	59.2	59.3	60.3	12	11	11	11
(5)	55.2	54.1	53.1	52.5	14	13	12	12
(6)	51.0	52.3	52.1	53.1	16	15	14	13
(7)	50.2	48.3	49.2	48.2	15	15	14	14
(8)	40.5	41.3	38.2	37.2	17	16	14	14
(9)	41.2	39.2	37.2	36.1	21	20	18	16

\*(1) LAD- CATREG, \*(2) WLAD-CATREG, \*(3) Lasso-CATREG, \*(4) adaptive lasso-CATREG, \*(5) WLAD adoptive lasso- CATREG, \*(6) naïve elastic net- CATREG, \*(7) Adaptive naïve elastic net-CATREG, \*(8) weighted naïve elastic net-CATREG, \*(9) WLAD naïve adoptive elastic net-CATREG.

In Table (1), the WLAD-CATREG method has a poor performance especially at high level of correlation and lasso-CATREG has working poor at high level of outliers. Parallel, the weighted naïve elastic net-CATREG and WLAD naïve adoptive elastic net-CATREG method have a best performance. In all cases, elastic net methods are more improved than lasso methods and the weight was given to improve the estimate methods. When the level of correlation and the level of outliers are increase, the elastic estimators.

### 6. Conclusions

We presented a new estimator for categorical regression model (CATREG) which takes into account effect the outliers, leverage point and variable select problem. The new criterion, Absolute Cross-validation (ACV) which use Cross-validation at robust form, was used for trade-off between estimators. The result for simulation study showed that, the weighted naïve elastic net-CATREG and the WLAD naïve adoptive elastic net-CATREG make improvements for the estimate.

### REFERENCES

[1] Peter, V. (1993), “M-estimators in multiple regression with optimal scaling”, Netherlands organization for scientific research, RR-92-10.  
 [2] Giloni, A., Simonoff, J., Sengupta, B., (2006a), “Robust weighted LAD regression”, Computational Statistics & Data Analysis, Vol.50.  
 [3] Giloni, A., Sengupta, B., Simonoff, J., (2006b), “A mathematical programming approach for improving the robustness of least sum of absolute deviations regression”,

Wiley InferScience doi, Vol.10.  
 [4] Olcay A., (2011), "Weighted LAD-LASSO method for robust parameter estimation and variable selection in regression", Computational Statistics & Data Analysis, Vo.56.  
 [5] Tibshirani, R. (1996), “Regression shrinkage and selection Via the lasso”, Journal of Royal Statistical Society B, Vo.58.  
 [6] Fan, J., and Li, R., (2001), “Variable selection via nonconcave penalized likelihood and its oracle properties”, Journal of the American Statistical Association, Vo. 96.  
 [7] Zou, H., (2006), “The adaptive Lasso and its oracle properties”, Journal of the American Statistical Association”, Vo. 101.  
 [8] Hansheng W., Chenlei L., (2006), "A note on adaptive group lasso", Computational Statistics and Data Analysis, Vo.52.  
 [9] Xu J., Yi Z., (2010), "Simultaneous estimation and variable selection in median regression using Lasso-type penalty", Annals of the institute of Statistical Mathematics, Vo.62.  
 [10] Zou, H. and Hastie, T., (2005), “Regularization and variable selection via the elastic net”. Journal of the Royal Statistical Society, Series B, Vo.67.  
 [11] Samiran G., (2007), "Adaptive Elastic Net: An Improvement of Elastic Net to achieve Oracle Properties", <http://www.math.iupui.edu/research/preprint/2007/p07-01.pdf>.  
 [12] Anita J., Van d., (2007), "Prediction Accuracy and Stability of Regression with Optimal Scaling Transformations", [https://www.researchgate.net/publication/28648961\\_Prediction\\_Accuracy\\_and\\_Stability\\_of\\_Regressionwith\\_Optimal\\_Scaling\\_Transformations](https://www.researchgate.net/publication/28648961_Prediction_Accuracy_and_Stability_of_Regressionwith_Optimal_Scaling_Transformations).  
 [13] Rousseeuw p. and Leroy A., (1987), "Robust Regression and Outlier Detection", New York: John Wiley.  
 [14] Ellis, S., Morgenthaler, S., (1992), "Leverage and break down in L1 regression", Journal of the American Statistical Association Vo. 87.

- [15] Chatterjee, S. and Hadi, A.S. (1988), *Sensitivity Analysis in Linear Regression*, Wiley, NewYork.
- [16] Rousseeuw p. and Hubert M., (1997), "Recent development in PROGRESS", <http://win-www.uia.ac.be/u/statis>.
- [17] Zou H., Zhang H. (2009), "On the adaptive elastic net with a diverging number of parameters", *Institute of Mathematical Statistics*, Vo.37, No.4.
- [18] Wang, H., Li, G., Jiang, G., (2007), "Robust regression shrinkage and consistent variable selection through the LAD-Lasso", *Journal of Business & Economic Statistics*, Vo.25.
- [19] Hong, D., Zhang, F., (2010), "Weighted elastic net model for mass spectrometry imaging processing. Mathematical modeling in the medical sciences", Vo.5, No.3.
- [20] Li, J.-T. , Jia, Y.-M., (2010), "An Improved Elastic Net for Cancer Classification and Gene Selection", *Acta Automatica Sinica*, Vo.36, No.7.
- [21] Wang, H., Leng, C., (2007), "Unified lasso estimation by least squares approximation", *Journal of the American Statistical Association*, Vo.102.
- [22] Pollard, D. (1991), "Asymptotics for least absolute deviation regression estimators", *Econometric Theory*, Vo.7.
- [23] Knight' K., (1998), "Limiting Distributions for L1 Regression Estimators under General Conditions. *Annals of Statistics*", Vo. 26, No.2.