

# DIVERGENCE MEASURES AND LOGISTIC REGRESSION MODELS

Leandro Pardo

**Abstract.** In this paper we present a review of some results about inference based on  $\phi$ -divergence measures, under assumptions of logistic regression model . The minimum  $\phi$ -divergence estimator, which is seen to be a generalization of the maximum likelihood estimator is considered. This estimator is used in a  $\phi$ -divergence measure which is the basis of new statistics for solving some important problems regarding logistic regression models: fitting the logistic regression model, residuals and dimensional reduction. Finally, an extension is presented when we consider a multinomial response instead a binary response.

**Keywords:** Chi-square distribution, Logistic regression models, Minimum  $\phi$ -divergence estimator, Goodness-of fit tests

**AMS classification:** 62B10, 62J15

## §1. Introduction

Let  $Y_i, i = 1, \dots, I$ , be independent binomial random variables with parameters  $n_i$  and  $\mathbf{p}_i, i = 1, \dots, I$ , we denote by  $n_{i1}$  the number of "succes". We assume that  $\mathbf{p}_i \equiv \pi(\mathbf{x}_i^T \beta)$ , where  $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ik})^T$  is a vector of  $k + 1$  explanatory variables and the dependence of them with  $\mathbf{p}_i$  is given by,

$$\pi(\mathbf{x}_i^T \beta) = \exp(\mathbf{x}_i \beta) / (1 + \exp(\mathbf{x}_i \beta)),$$

where  $\beta = (\beta_0, \beta_1, \dots, \beta_k)^T$  is a  $k + 1$  vector of unknown regression coefficients. The parameter space is given by

$$\Theta = \{(\beta_0, \beta_1, \dots, \beta_k) : \beta_i \in (-\infty, \infty), i = 0, \dots, k\} \tag{1}$$

and the loglikelihood function for  $\beta$  by

$$\begin{aligned} \log L(\beta) &= \log \prod_{i=1}^I \binom{n_i}{n_{i1}} \pi(\mathbf{x}_i^T \beta)^{n_{i1}} (1 - \pi(\mathbf{x}_i^T \beta))^{n_i - n_{i1}} \\ &= c - \left\{ \sum_{i=1}^I \left( \frac{n_{i1}}{N} \log \frac{n_{i1}}{\pi(\mathbf{x}_i^T \beta) n_i} + \frac{n_i - n_{i1}}{N} \log \frac{n_i - n_{i1}}{(1 - \pi(\mathbf{x}_i^T \beta)) n_i} \right) \right\}, \end{aligned}$$

where  $N = n_1 + \dots + n_I$ .

We consider the two following probability vectors

$$\widehat{\mathbf{p}} = \left( \frac{n_{11}}{N}, \frac{n_{12}}{N}, \frac{n_{21}}{N}, \frac{n_{22}}{N}, \dots, \frac{n_{I1}}{N}, \frac{n_{I2}}{N} \right)^T$$

and

$$\mathbf{p}(\beta) \equiv \left( \pi(\mathbf{x}_1^T \beta) \frac{n_1}{N}, (1 - \pi(\mathbf{x}_1^T \beta)) \frac{n_1}{N}, \dots, \pi(\mathbf{x}_I^T \beta) \frac{n_I}{N}, (1 - \pi(\mathbf{x}_I^T \beta)) \frac{n_I}{N} \right)^T.$$

It is immediate to get

$$\log L(\beta) = c - D_{Kull}(\widehat{\mathbf{p}}, \mathbf{p}(\beta)),$$

where by  $D_{Kull}(\widehat{\mathbf{p}}, \mathbf{p}(\beta))$  we are denoting the Kullback-Leibler divergence between the probability vectors  $\widehat{\mathbf{p}}$  and  $\mathbf{p}(\beta)$ . Then the maximum likelihood estimator of  $\beta$  can be defined as

$$\widehat{\beta} = \arg \min_{\beta_0, \beta_1, \dots, \beta_k} D_{Kullback}(\widehat{\mathbf{p}}, \mathbf{p}(\beta)). \quad (2)$$

## §2. Minimum $\phi$ -divergence estimator

The Kullback-Leibler divergence measure is a particular case of the  $\phi$ -divergence measures defined by Csiszár (1963) and Ali and Silvey (1966), as

$$D_\phi(\widehat{\mathbf{p}}, \mathbf{p}(\beta)) = \left\{ \sum_{i=1}^I \left( \frac{\pi(\mathbf{x}_i^T \beta) n_i}{N} \phi \left( \frac{n_{i1}}{\pi(\mathbf{x}_i^T \beta) n_i} \right) + \frac{(1 - \pi(\mathbf{x}_i^T \beta)) n_i}{N} \phi \left( \frac{n_i - n_{i1}}{(1 - \pi(\mathbf{x}_i^T \beta)) n_i} \right) \right) \right\}, \quad (3)$$

where  $\phi \in \Phi^*$ , being  $\Phi^*$  the class of all convex functions  $\phi(x)$ ,  $x > 0$ , such that at  $x = 1$ ,  $\phi(1) = 0$ ,  $\phi''(1) > 0$ , and at  $x = 0$ ,  $0\phi(0/0) = 0$  and  $0\phi(p/0) = \lim_{u \rightarrow \infty} \phi(u)/u$ . For every  $\phi \in \Phi^*$  that is differentiable at  $x = 1$ , the function  $\psi(x) \equiv \phi(x) - \phi'(1)(x - 1)$  also belongs to  $\Phi^*$ . Then we have  $D_\psi(\widehat{\mathbf{p}}, \mathbf{p}(\beta)) = D_\phi(\widehat{\mathbf{p}}, \mathbf{p}(\beta))$ , and  $\psi$  has the additional property that  $\psi'(1) = 0$ . Because the two divergence measures are equivalent, we can consider the set  $\Phi^*$  to be equivalent to the set  $\Phi \equiv \Phi^* \cap \{\phi : \phi'(1) = 0\}$ . In what follows, we give our theoretical results for  $\phi \in \Phi$ , but we often apply them to choices of functions in  $\Phi^*$ .

Based on (2) and on the definition of  $D_\phi(\widehat{\mathbf{p}}, \mathbf{p}(\beta))$ , in Pardo et al (2003a), it was defined and studied the minimum  $\phi$ -divergence estimator of  $\beta$ . This estimator is defined by

$$\widehat{\beta}^\phi = \arg \min_{\alpha, \beta_1, \dots, \beta_k} D_\phi(\widehat{\mathbf{p}}, \mathbf{p}(\beta)).$$

For  $\phi(x) = x \log x - x + 1$  we obtain the maximum likelihood estimator and for  $\phi(x) = \frac{1}{2}(x - 1)^2$  the minimum chi-squared estimator.

If we denote by  $\mathbf{X}$  the  $I \times (k + 1)$  matrix with rows  $\mathbf{x}_i$ ,  $i = 1, \dots, I$  and we assume that  $\text{rank}(\mathbf{X}) = k + 1$ , in the cited paper of Pardo et al (2003a) it was established that under the assumption that  $\phi$  is twice continuously differentiable in a neighborhood of 1 that

$$\begin{aligned} \widehat{\beta}^\phi &= \beta^0 + \left( \mathbf{X}^T \text{Diag} \left( \left( \frac{n_i}{N} \pi(\mathbf{x}_i^T \beta^0) (1 - \pi(\mathbf{x}_i^T \beta^0)) \right)_{i=1, \dots, I} \right) \mathbf{X} \right)^{-1} \mathbf{X}^T \\ &\times \text{Diag} \left( \left( \mathbf{C}_i^T \right)_{i=1, \dots, I} \right) \text{Diag} \left( \mathbf{p}(\beta^0)^{-1/2} \right) (\widehat{\mathbf{p}} - \mathbf{p}(\beta^0)) + o(\|\widehat{\mathbf{p}} - \mathbf{p}(\beta^0)\|), \end{aligned}$$

where

$$C_i = \left( \left( \frac{n_i}{N} \right) \pi(\mathbf{x}_i^T \beta^0) (1 - \pi(\mathbf{x}_i^T \beta^0)) \right)^{1/2} \begin{pmatrix} (1 - \pi(\mathbf{x}_i^T \beta^0))^{1/2} \\ -\pi(\mathbf{x}_i^T \beta^0) \end{pmatrix}, i = 1, \dots, I$$

and also, under the assumption that  $n_i/n \rightarrow \lambda_i$  when  $n_i \rightarrow \infty$ , we have

$$\sqrt{N} (\widehat{\beta}^\phi - \beta^0) \xrightarrow[N \rightarrow \infty]{L} N \left( \mathbf{0}, \left( \mathbf{X}^T \text{Diag} (\lambda_i \pi(\mathbf{x}_i^T \beta^0) (1 - \pi(\mathbf{x}_i^T \beta^0)))_{i=1, \dots, I} \mathbf{X} \right)^{-1} \right) \quad (4)$$

being  $\lambda_i = \lim_{N \rightarrow \infty} n_i/N$ .

From (4) it is clear that

$$\widehat{\text{Cov}}(\widehat{\beta}^\phi) \approx \frac{1}{N} \left( \mathbf{X}^T \text{Diag} (\lambda_i \pi(\mathbf{x}_i^T \widehat{\beta}^\phi) (1 - \pi(\mathbf{x}_i^T \widehat{\beta}^\phi)))_{i=1, \dots, I} \mathbf{X} \right)^{-1}.$$

### §3. Fitting the Logistic Regression Model: Residuals

We denote by  $\widehat{\beta}_0, \widehat{\beta}_1, \dots, \widehat{\beta}_k$  the maximum Likelihood estimators of  $\beta_0, \beta_1, \dots, \beta_k$ . If we estimate  $\pi(\mathbf{x}_i^T \beta)$  by the maximum likelihood estimator  $\pi(\mathbf{x}_i^T \widehat{\beta})$ , the classical Pearson chi-squared statistic is given by

$$X^2 = \sum_{i=1}^I \frac{(n_{i1} - n_i \pi(\mathbf{x}_i^T \widehat{\beta}))^2}{n_i \pi(\mathbf{x}_i^T \widehat{\beta}) (1 - \pi(\mathbf{x}_i^T \widehat{\beta}))}$$

and the likelihood ratio test statistic by

$$D^2 = \sum_{i=1}^I 2 \left\{ n_{i1} \log \frac{n_{i1}}{n_i \pi(\mathbf{x}_i^T \widehat{\beta})} + (n_i - n_{i1}) \log \frac{n_i - n_{i1}}{n_i (1 - \pi(\mathbf{x}_i^T \widehat{\beta}))} \right\}.$$

If now we consider the minimum  $\phi_2$ -divergence estimators  $\widehat{\beta}_0^{\phi_2}, \widehat{\beta}_1^{\phi_2}, \dots, \widehat{\beta}_k^{\phi_2}$  of parameters  $\beta_0, \beta_1, \dots, \beta_k$ , instead of the maximum likelihood estimators, we can estimate  $\pi(\mathbf{x}_i^T \beta)$  by  $\pi(\mathbf{x}_i^T \widehat{\beta}^{\phi_2})$ . Pardo et al (2003b) have considered and studied the  $\phi_1$ -divergence test statistic based on the minimum  $\phi_2$ -divergence estimator  $\widehat{\beta}^{\phi_2}$ , given by

$$\begin{aligned} T_{\phi_1, \phi_2} &= \frac{2N}{\phi_1''(1)} D_{\phi_1}(\widehat{\mathbf{p}}, \mathbf{p}(\widehat{\beta}^{\phi_2})) \\ &= \frac{2}{\phi_1''(1)} \sum_{i=1}^I n_i \left\{ \pi(\mathbf{x}_i^T \widehat{\beta}^{\phi_2}) \phi \left( \frac{n_{i1}}{\pi(\mathbf{x}_i^T \widehat{\beta}^{\phi_2}) n_i} \right) \right. \\ &\quad \left. + \left( 1 - \pi(\mathbf{x}_i^T \widehat{\beta}^{\phi_2}) \right) \phi \left( \frac{n_{i2}}{(1 - \pi(\mathbf{x}_i^T \widehat{\beta}^{\phi_2})) n_i} \right) \right\}. \end{aligned} \quad (5)$$

It is interesting to observe that for

$$\phi_2(x) = x \log x - x + 1 \quad \text{and} \quad \phi_1(x) = \frac{1}{2} (x - 1)^2,$$

we obtain that  $T_{\phi_1, \phi_2} \equiv X^2$  and for

$$\phi_2(x) = x \log x - x + 1 \quad \text{and} \quad \phi_1(x) = x \log x - x + 1,$$

we get  $T_{\phi_1, \phi_2} \equiv D^2$ . The following theorem presents the asymptotic distribution of the family of test statistics  $T_{\phi_1, \phi_2}$ .

**Theorem 1.** *Suppose that the data  $Y_i, i = 1, \dots, I$  are binomially distributed with parameters  $n_i$  and  $\pi(\mathbf{x}_i^T \beta)$ . Choose functions  $\phi_1$  and  $\phi_2 \in \Phi$  and twice continuously differentiable in a neighborhood of 1. Under the hypothesis  $\mathbf{p} = \mathbf{p}(\beta)$  and assuming that  $n_i/n \rightarrow \lambda_i > 0$  when  $n_i \rightarrow \infty$ , the test statistic  $T_{\phi_1, \phi_2}$  has a chi-squared distribution with  $I - (k + 1)$  degrees of freedom.*

Based on this theorem, if the sample sizes are large enough, one can use the asymptotic quantile  $\chi_{I-(k+1), 1-\alpha}^2$ , defined by the equation  $P(\chi_{M-1}^2 \leq \chi_{M-1, 1-\alpha}^2) = 1 - \alpha$ , to propose the decision rule:

$$\text{“Reject } H_{Null} : \mathbf{p} = \mathbf{p}(\beta) \text{ if } T_{\phi_1, \phi_2} > \chi_{I-(k+1), 1-\alpha}^2 \text{”}.$$

The statistics  $X^2, D^2$  and  $T_{\phi_1, \phi_2}$  provide a single number which summarizes the agreement of observed and fitted values. The advantage (as well as the disadvantage) of these statistics is that a single number is used to summarize considerable information. Additional diagnostic analyses are necessary to describe the nature of one lack of fit. Residual comparing observed and fitted counts are useful for this purpose. From a classical point of view we have the residual based on the Pearson's test statistic  $X^2$ , given by

$$e_i = \frac{n_{i1} - n_i \pi(\mathbf{x}_i^T \hat{\beta})}{\sqrt{n_i \pi(\mathbf{x}_i^T \hat{\beta}) (1 - \pi(\mathbf{x}_i^T \hat{\beta}))}},$$

and residual based in the likelihood ratio test given by

$$d_i = \text{sig} \left( n_{i1} - n_i \pi(\mathbf{x}_i^T \hat{\beta}) \right) \sqrt{2 \left\{ n_{i1} \log \frac{n_{i1}}{n_i \pi(\mathbf{x}_i^T \hat{\beta})} + (n_i - n_{i1}) \log \frac{n_i - n_{i1}}{n_i (1 - \pi(\mathbf{x}_i^T \hat{\beta}))} \right\}}.$$

In the same way we can define  $\phi_1$ -residuals based on minimum  $\phi_2$ -divergence estimator by

$$\begin{aligned} c_i^{\phi_1, \phi_2} = & \text{sig} \left( n_{i1} - n_i \pi(\mathbf{x}_i^T \hat{\beta}^{\phi_2}) \right) \sqrt{\frac{2n_i}{\phi_1''(1)}} \left\{ \pi(\mathbf{x}_i^T \hat{\beta}^{\phi_2}) \phi_1 \left( \frac{n_{i1}}{\pi(\mathbf{x}_i^T \hat{\beta}^{\phi_2}) n_i} \right) \right. \\ & \left. + \left( 1 - \pi(\mathbf{x}_i^T \hat{\beta}^{\phi_2}) \right) \phi_1 \left( \frac{n_{i2}}{(1 - \pi(\mathbf{x}_i^T \hat{\beta}^{\phi_2})) n_i} \right) \right\}^{1/2}. \end{aligned}$$

In the following theorem we present its asymptotic distribution.

**Theorem 2.** *Suppose that the data  $Y_i, i = 1, \dots, I$  are binomially distributed with parameters  $n_i$  and  $\pi(\mathbf{x}_i^T \beta)$ . Choose functions  $\phi_1$  and  $\phi_2 \in \Phi$  twice continuously differentiable in a neighborhood of 1. Assuming that  $n_i/n \rightarrow \lambda_i > 0$  when  $n_i \rightarrow \infty$ , we have*

$$c_j^{\phi_1, \phi_2} \xrightarrow[N \rightarrow \infty]{L} \mathcal{N}(0, \tau_j^2)$$

where

$$\tau_j^2 = 1 - \lambda_j \pi(\mathbf{x}_j^T \beta^0) (1 - \pi(\mathbf{x}_j^T \beta^0)) \mathbf{x}_j \mathbf{I}_F(\beta^0)^{-1} \mathbf{x}_j^T$$

being  $\mathbf{I}_F(\beta^0) = \mathbf{X}^T \text{Diag}(\lambda_i \pi(\mathbf{x}_i^T \beta^0) (1 - \pi(\mathbf{x}_i^T \beta^0)))_{i=1, \dots, I} \mathbf{X}$  the Fisher Information matrix associated with the Logistic Regression Model.

This result is important because we can define the  $\phi_1$ -standardized residuals based on minimum  $\phi_2$ -divergence estimator by

$$(c_j^{\phi_1, \phi_2})^* = \frac{c_j^{\phi_1, \phi_2}}{\widehat{\tau}_j^{\phi_2}}$$

being

$$\left(\widehat{\tau}_j^{\phi_2}\right)^2 = 1 - \frac{n_j}{N} \pi(\mathbf{x}_j^T \widehat{\beta}^{\phi_2}) \left(1 - \pi(\mathbf{x}_j^T \widehat{\beta}^{\phi_2})\right) \mathbf{x}_j \mathbf{I}_F(\widehat{\beta}^{\phi_2})^{-1} \mathbf{x}_j^T.$$

### §4. Dimensionality reduction

By dimensionality reduction we understand the procedure to determine if the independent variables in the model are "significantly" related to the outcome variable. To choose this logistic regression model we use a backward deletion procedure. Any stepwise procedure for deletion of variables from a model is based on a statistical algorithm which checks for the "importance" of a variable is defined in terms of a measure of the statistical significance of the coefficient for the variable, i.e., we must carry out the hypothesis test

$$H_{Null} : \beta_j = 0 \text{ against } H_{Alt} : \beta_j \neq 0 \quad (j = 1, \dots, k) \tag{6}$$

In the first stage, we consider  $j = 1, \dots, k$ , i.e., the logistic regression model with all available explanatory variables. Then, we delete the explanatory variable associated with the regression parameter,  $\beta_{k_1}$ , if the associated  $p$ -value for testing, (6) with  $j = k_1, p_{k_1}$ , is the highest. In the second stage, we rename the parameters  $\beta_{k_1}, \dots, \beta_k$  as  $\beta_{k_1+1}, \dots, \beta_{k-1}$ , respectively. Therefore, the logistic regression model with all the explanatory variables except the corresponding to the parameter  $\beta_{k_1}$  is considered. So, after testing (6) for  $j = 1, \dots, k - 1$ , we delete the explanatory variable with the highest associated  $p$ -value and so on. Finally, we stop the procedure when the maximum  $p$ -value associated with the logistic regression model with the remaining explanatory variables is sufficiently small.

In Pardo et al. (2003b) in order to solve the problem considered in (6) the following families of tests statistics were considered

$$S_{\phi_1, \phi_2}^{\beta_j, t} = \frac{2N}{\phi_1''(1)} D_{\phi_1} \left( \mathbf{p}(\widehat{\beta}^{\phi_2, t}), \mathbf{p}(\widehat{\beta}^{\phi_2, t}) \right),$$

where  $\widehat{\beta}^{\phi_2,t}$  is the minimum  $\phi_2$ -divergence estimator of  $\beta_1, \dots, \beta_{k+1-t}$  and  ${}^j\widehat{\beta}^{\phi_2}$  is the minimum  $\phi_2$ -divergence estimator of  $(\beta_0, \dots, \beta_{j-1}, 0, \beta_{j+1}, \dots, \beta_{k+1-t})$  in the stage  $t$ , and

$$T_{\phi_1, \phi_2}^{\beta_j, t} = \frac{2N}{\phi_1''(1)} \left( D_{\phi_1} \left( \widehat{\mathbf{p}}, \mathbf{p} \left( \widehat{\beta}^{\phi_2, t} \right) \right) - D_{\phi_1} \left( \widehat{\mathbf{p}}, \mathbf{p} \left( {}^j\widehat{\beta}^{\phi_2, t} \right) \right) \right).$$

In the following theorem we present the asymptotic distribution of the tests statistic  $S_{\phi_1, \phi_2}^{\beta_j, t}$  and  $T_{\phi_1, \phi_2}^{\beta_j, t}$ .

**Theorem 3.** *Suppose that the data  $Y_i, i = 1, \dots, I$  are binomially distributed with parameters  $n_i$  and  $\pi(\mathbf{x}_i^T \beta)$ . Choose functions  $\phi_1$  and  $\phi_2 \in \Phi$  twice continuously differentiable in a neighborhood of 1. Then for testing*

$$H_{Null} : \beta_j = 0 \text{ versus } H_{Alt} : \beta_j \neq 0,$$

*the test statistics  $S_{\phi_1, \phi_2}^{\beta_j, t}$  and  $T_{\phi_1, \phi_2}^{\beta_j, t}$  are asymptotically distributed as a chi-squared distribution with 1 degree of freedom, under the assumption that  $n_i/n \rightarrow \lambda_i > 0$  when  $n_i \rightarrow \infty$ .*

Based on this theorem, if the sample sizes are large enough, one can use the asymptotic quantile  $\chi_{1,1-\alpha}^2$ , defined by the equation  $P(\chi_{M-1}^2 \leq \chi_{M-1,1-\alpha}^2) = 1 - \alpha$ , to propose the decision rule:

$$\text{“Reject } H_{Null} \text{ if } S_{\phi_1, \phi_2}^{\beta_j, t} > \chi_{M-1,1-\alpha}^2 \text{ (or } T_{\phi_1, \phi_2}^{\beta_j, t} > \chi_{1,1-\alpha}^2 \text{)”}.$$

### §5. Polytomous Logistic Regression Model

We consider a response random variable  $Y$  belonging to one of the  $J$  distinct categories  $C_1, \dots, C_J$ , which is observed together with  $p + 1$  explanatory variables  $\mathbf{x}^T = (1, x_1, \dots, x_p) \in \mathbb{R}^{p+1}$ . For convenience  $x_0 = 1$ . Let  $\pi_j(\mathbf{x}) = P(Y \in C_j | \mathbf{x}), j = 1, \dots, J$ , denote the probability that the observation of the random variable  $Y$  belongs to the category  $C_j, j = 1, \dots, J$ , when the explanatory variable is  $\mathbf{x}^T$ . More specifically suppose here, that the dependence between  $Y$  and  $\mathbf{x}^T$  can be modeled by using the logistic assumption

$$\pi_j(\mathbf{x}) \equiv \pi(\mathbf{x}^T \beta_j) = \exp(\mathbf{x}^T \beta_j) / \sum_{l=1}^J \exp(\mathbf{x}^T \beta_l), \quad j = 1, \dots, J, \tag{7}$$

where  $\beta_j^T = (\beta_{0j}, \dots, \beta_{pj}), j = 1, \dots, J - 1$ , is a vector of unknown parameters and  $\beta_J^T = (0, \dots, 0)$ , for convenience. The vector  $\beta^T = (\beta_1^T, \dots, \beta_{J-1}^T)$  is  $\nu$ -dimensional with  $\nu = (J - 1)(p + 1)$ . The model described in (7) is the classical *Polytomous Logistic Regression Model* (PLRM) or *Multinomial Logistic Regression Model*. For more details about this model see Amemiya (1981), Anderson (1972, 1982, 1984), Lesaffre (1986), Lesaffre and Albert (1986, 1989), Mantel (1966), Theil (1969), McCullag (1980) Agresti (2002), Engel (1988) and references there in. In the following we shall denote by

$$\Theta = \{ \beta_j^T = (\beta_{0j}, \dots, \beta_{pj}), j = 1, \dots, J - 1 : \beta_{st} \in \mathbb{R}, s = 0, \dots, p; t = 1, \dots, J - 1 \}$$

We assume that  $N$  different values of the vector of explanatory variables,

$$\mathbf{x}_i^T = (x_{i0}, x_{i1}, \dots, x_{ip}), \quad i = 1, \dots, N,$$

are available. Let  $n(\mathbf{x}_i)$  be the number of observations considered when the explanatory variable  $\mathbf{x}^T$  has the value  $\mathbf{x}_i^T$ , so that if  $\mathbf{x}^T$  is fixed at  $\mathbf{x}_i^T$  we have a multinomial distribution with parameters

$$(n(\mathbf{x}_i); \pi(\mathbf{x}_i^T \beta_1), \dots, \pi(\mathbf{x}_i^T \beta_J)).$$

Vectors of probabilities are denoted by  $\pi(\mathbf{x}_i) = (\pi(\mathbf{x}_i^T \beta_1), \dots, \pi(\mathbf{x}_i^T \beta_J))^T$  and total sample size by  $n = n(\mathbf{x}_1) + \dots + n(\mathbf{x}_N)$ . Given the explanatory variable  $\mathbf{x}_i^T$ , we denote the number of observations in the class  $C_s$  by  $y_{si}$ . It is clear that  $n(\mathbf{x}_i) = \sum_{s=1}^J y_{si}$ . To estimate  $\beta_{js}$  ( $j = 0, \dots, p; s = 1, \dots, J - 1$ ) we maximize the loglikelihood function

$$\log \prod_{i=1}^N \prod_{l=1}^J \frac{n(\mathbf{x}_i)!}{y_{1i}! \times \dots \times y_{Ji}!} \pi(\mathbf{x}_i^T \beta_l)^{y_{li}} \approx \log \prod_{i=1}^N \prod_{l=1}^J \pi(\mathbf{x}_i^T \beta_l)^{y_{li}} \equiv L(\beta)$$

with  $y_{Ji} = n(\mathbf{x}_i) - \sum_{s=1}^{J-1} y_{si}$ .

It is not difficult to establish that

$$\left( \frac{\partial^2 L(\beta)}{\partial \beta^2} \right)_{(p+1)(J-1) \times (p+1)(J-1)} = -n \mathbf{X}^T \mathbf{V}_n(\beta) \mathbf{X},$$

where

$$\mathbf{X}^T = (\mathbf{X}_1^T, \dots, \mathbf{X}_N^T)_{(p+1)(J-1) \times (J-1)N},$$

being

$$\mathbf{X}_i = \begin{pmatrix} \mathbf{x}_i^T & \mathbf{0}^T & \dots & \mathbf{0}^T \\ \mathbf{0}^T & \mathbf{x}_i^T & \dots & \mathbf{0}^T \\ \cdot & \cdot & \cdot & \cdot \\ \mathbf{0}^T & \mathbf{0}^T & \dots & \mathbf{x}_i^T \end{pmatrix}_{(J-1) \times (J-1)(p+1)} \tag{8}$$

and the matrix  $\mathbf{V}_n(\beta)$  is defined by

$$\mathbf{V}_n(\beta) = \text{diag} \left( \frac{n(\mathbf{x}_1)}{n} \mathbf{V}_1(\beta), \dots, \frac{n(\mathbf{x}_N)}{n} \mathbf{V}_N(\beta) \right)_{N(J-1) \times N(J-1)}$$

with

$$\mathbf{V}_i(\beta)_{(J-1) \times (J-1)} = (\pi(\mathbf{x}_i^T \beta_s) (\delta_{st} - \pi(\mathbf{x}_i^T \beta_t)))_{s,t=1,\dots,J-1}, \quad i = 1, \dots, N,$$

and  $\delta_{st}$  is the Kronecker delta.

The Fisher information matrix is given by

$$\mathbf{I}_{F,n}(\beta) = \mathbf{X}^T \mathbf{V}_n(\beta) \mathbf{X} = \sum_{j=1}^N \frac{n(\mathbf{x}_j)}{n} \mathbf{X}_j^T \mathbf{V}_j(\beta) \mathbf{X}_j$$

and

$$\sqrt{n} (\hat{\beta} - \beta_0) \xrightarrow[n \rightarrow \infty]{L} N(\mathbf{0}, \mathbf{I}_{F,\lambda}(\beta_0)),$$

where  $\hat{\beta}$  is the maximum likelihood estimator of  $\beta$  and  $\beta_0$  is the true value of the parameter  $\beta$ ,

$$\mathbf{I}_{F,\lambda}(\beta_0) = \sum_{j=1}^N \lambda_j \mathbf{X}_j^T \mathbf{V}_j(\beta_0) \mathbf{X}_j \quad \text{and} \quad \lambda_j = \lim_{n \rightarrow \infty} \frac{n(\mathbf{x}_j)}{n}, \quad j = 1, \dots, N.$$

Let us introduce the two following probability vectors,

$$\widehat{\mathbf{p}} = \left( \frac{y_{11}}{n}, \dots, \frac{y_{J1}}{n}, \frac{y_{12}}{n}, \dots, \frac{y_{J2}}{n}, \dots, \frac{y_{1N}}{n}, \dots, \frac{y_{JN}}{n} \right)^T,$$

and

$$\begin{aligned} \mathbf{p}(\beta) &= \left( \frac{n(\mathbf{x}_1)}{n} (\pi(\mathbf{x}_1^T \beta_1), \dots, \pi(\mathbf{x}_1^T \beta_J)), \dots, \frac{n(\mathbf{x}_N)}{n} (\pi(\mathbf{x}_N^T \beta_1), \dots, \pi(\mathbf{x}_N^T \beta_J)) \right)^T \\ &= \left( \frac{n(\mathbf{x}_1)}{n} \pi(\mathbf{x}_1)^T, \dots, \frac{n(\mathbf{x}_N)}{n} \pi(\mathbf{x}_N)^T \right)^T. \end{aligned} \quad (9)$$

The Kullback-Leibler divergence measure between the probability vectors  $\widehat{\mathbf{p}}$  and  $\mathbf{p}(\beta)$  is

$$\begin{aligned} D_{Kullback}(\widehat{\mathbf{p}}, \mathbf{p}(\beta)) &= \sum_{l=1}^J \sum_{i=1}^N \frac{y_{li}}{n} \log \frac{\frac{y_{li}}{n}}{\pi(\mathbf{x}_i^T \beta_l) \frac{n(\mathbf{x}_i)}{n}} \\ &= kte - \frac{1}{n} \log \prod_{i=1}^N \prod_{l=1}^J \pi(\mathbf{x}_i^T \beta_l)^{y_{li}} \\ &\approx -L(\beta). \end{aligned}$$

Then the MLE of parameter  $\beta$  can be equivalently defined by the condition

$$\widehat{\beta} = \arg \min D_{Kullback}(\widehat{\mathbf{p}}, \mathbf{p}(\beta)),$$

and its extension to the minimum  $\phi$ -divergence estimator, is given by

$$\widehat{\beta}^\phi \equiv \arg \min_{\beta_{01}, \dots, \beta_{pJ-1}} D_\phi(\widehat{\mathbf{p}}, \mathbf{p}(\beta)). \quad (10)$$

In relation to the asymptotic properties of the minimum  $\phi$ -divergence estimator we have that the most important result is

$$\sqrt{n} (\widehat{\beta}^\phi - \beta_0) \xrightarrow[n \rightarrow \infty]{L} N(\mathbf{0}, \mathbf{I}_{F,\lambda}(\beta_0)),$$

and then

$$Cov(\widehat{\beta}^\phi) \approx \frac{1}{N} (\mathbf{X}^T \mathbf{V}_n(\widehat{\beta}^\phi) \mathbf{X})^{-1}.$$

After estimating the unknown parameters, we would like to know how effective the model we have is in describing the outcome variable. This is referred to as goodness-of-fit. We will conclude that the model fits if (a) summary measures of the distance between the observed sample values and the values predicted by the model are small and (b) the contribution of each pair (observed, predicted) to these summary measures is unsystematic and is small relative to the error structure of the model. Thus, a complete assessment of the fitted model will involve both (a) computation and evaluation of overall measures of fit, (b) examination of the individual components of these measures. It is possible to present families of statistics based on  $\phi$ -divergence measures to solve (a) in a similar way to the families of test statistics given in (5). Additional diagnosis analyses are necessary to describe the nature of one lack of fit. A family of residuals based on  $\phi$ -divergences that is a generalization of the classical residuals it is possible to define on the basis of  $\phi$ -divergence measures.



## References

- [1] ALI, S.M. AND SILVEY, S.D. A general class of coefficient of divergence of one distribution from another. *Journal of Royal Statistical Society, Series B*, 286 (1966), 131–142.
- [2] AGRESTI, A. *Categorical Data Analysis* (Second Edition), John Wiley & Sons, 2002.
- [3] AMEMIYA, T. Qualitative response models: a survey. *Journal of Economic Literature*, 19 (1981), 1483-1536.
- [4] ANDERSON, J. A. Separate sample logistic discrimination. *Biometrika*, 59 (1972), 19-35.
- [5] ANDERSON, J. A. Logistic discrimination. In *Handbook of Statistics*, P. R. Krishnaiah and L. N. Kanal, Eds., North-Holland Publ. Comp. , 1982, 169-191.
- [6] ANDERSON, J. A. Regression and ordered categorical variables. *Journal of Royal Statistical Society, Series B*, 46 (1984), 1-30.
- [7] CSISZAR, I.. Eine Informationstheoretische Ungleichung und ihre Anwendung auf den Beweis der Ergodizität von Markoffschen Ketten. *Publications of the Mathematical Institute of Hungarian Academy of Sciences, Series A*, 8 (1963), 85–108.
- [8] ENGEL, J. Polytomous logistic regression. *Statistica Neerlandica*, 42 (1988), 233-252.
- [9] KULLBACK, S. Kullback information. In *Encyclopedia of Statistical Sciences, Volume 4*, editors S. Kotz and N. L. Johnson, John Wiley & Sons, New York, 1985, pp. 421-425.
- [10] LESAFFRE, E. Logistic discrimination analysis with application in electrocardiography, *P.H. Dissertation*. University of Leuven, 1986.
- [11] LESAFFRE, E. AND ALBERT, A. Multiple-group Logistic Regression Diagnostic. *Applied Statistics*, 38 (1989), 425-440.
- [12] MCCULLAG, P. Regression Models for Ordinary Data. *Journal of the Royal Statistical Society*, B 42 (1980), 109-142.
- [13] MANTE, N. Models for complex contingency tables and polychotomous dosage response curves. *Biometrics*, 22 (1966), 83-95.
- [14] PARDO, J. A., PARDO, L. AND PARDO, M. C. Testing in logistic regression models based on  $\phi$ -divergences measures. Submitted, (2003)
- [15] PARDO, J. A., PARDO, L. AND PARDO, M. C. Minimum  $\phi$ -divergence estimator in logistic regression models. Submitted, (2003)

Leandro Pardo

Department of Statistics and O.R.,

Complutense University of Madrid.

telephone: 34-913944425 fax: 34-913944607

Leandro\_Pardo@Mat.ucm.es