

반파라메터최량판별함수에 대한 추정

림 창 호

경애하는 최고령도자 김정은동지께서는 다음과 같이 말씀하시였다.

《이미 일정한 토대가 있고 전망이 확고한 연구대상들에 힘을 넣어 세계패권을 쥐며 그 성과를 확대하는 방법으로 과학기술을 빨리 발전시켜야 합니다.》(《조선로동당 제7차대회에서 한 중앙위원회사업총화보고》 단행본 39페이지)

반파라메터판별함수에 의한 판별분석문제는 패턴인식을 비롯한 현실문제들에서 많이 제기된다.

선행연구에서는 회귀분석문제로서 부분선형회귀모형과 반파라메터회귀모형의 추론문제[7-9]와 파라메터판별함수에 의한 판별분석문제[3, 4, 6]를 고찰하였다.

또한 비파라메터판별함수에 의한 판별분석문제[1, 2]와 여러반파라메터모형들사이의 판별분석문제[5], 극값판별분석과 선형판별분석모형의 일반화로서 반파라메터선형판별분석모형[10, 11]을 논의하였다.

판별변량 x 에 관한 모집단 G_r ($r=1, 2$)의 밀도함수와 사전확률을 각각 $f_r^0(x)$, q_r^0 이라고 하자.

표본공간 R 를 공통부분이 없는 2개의 구역 R_1, R_2 로 나누고

$$x \in R_1 \Rightarrow x \in G_1, \quad x \in R_2 \Rightarrow x \in G_2$$

로 판별할 때 오판별확률

$$P = q_1^0 \int_{R_2} f_1^0(x) dx + q_2^0 \int_{R_1} f_2^0(x) dx$$

가 최소로 되는 최량판별규칙은

$$2F^0(x) > 1 \Rightarrow x \in G_1, \quad 2F^0(x) \leq 1 \Rightarrow x \in G_2 \quad (1)$$

이다. 여기서

$$F^0(x) = P\{G_1 | x\} = \frac{q_1^0 f_1^0(x)}{q_1^0 f_1^0(x) + q_2^0 f_2^0(x)}$$

이며 이때 $F^0(x)$ 를 최량판별함수라고 부른다.

모집단 G_r ($r=1, 2$)의 밀도함수 $f_r^0(x)$ ($r=1, 2$)와 사전확률 q_r^0 ($r=1, 2$)들이 미지라고 하면 최량판별규칙 (1)은 미지판별함수 $F(x) = q_1 f_1(x) / [q_1 f_1(x) + q_2 f_2(x)]$ 에 의한 판별규칙 $2F(x) > 1 \Rightarrow x \in G_1$, $2F(x) \leq 1 \Rightarrow x \in G_2$ 로 된다. 그리고 미지판별함수 $F(x)$ 는 미지파라메터 $\theta = (\theta_1, \dots, \theta_s)^T \in \Theta$ 와 미지함수 $\mathbf{g} = (g_1, \dots, g_q)^T \in H \equiv H(\mathbf{R}, \mathbf{R}^q)$ 를 포함하는 반파라메터판별함수 $F(x) = F(x; \theta, \mathbf{g})$ 로 된다.

이때 $F(x; \theta, \mathbf{g})$ 에 의한 판별규칙 $2F(x; \theta, \mathbf{g}) > 1 \Rightarrow x \in G_1$, $2F(x; \theta, \mathbf{g}) \leq 1 \Rightarrow x \in G_2$ 에서 오판별확률

$$P = q_1^0 \int_{R_2} f_1^0(x) dx + q_2^0 \int_{R_1} f_2^0(x) dx, \quad R_1 = \{x \mid 2F(x; \theta, \mathbf{g}) > 1\}, \quad R_2 = \{x \mid 2F(x; \theta, \mathbf{g}) \leq 1\}$$

이 최소로 되는 최량판별 규칙

$$2F(x; \theta^*, \mathbf{g}^*) > 1 \Rightarrow x \in G_1, \quad 2F(x; \theta^*, \mathbf{g}^*) \leq 1 \Rightarrow x \in G_2 \quad (2)$$

를 반파라메터최량판별 규칙, $F(x; \theta^*, \mathbf{g}^*)$ 을 반파라메터최량판별 함수라고 부른다.

논문에서는 표본자료 x_1, \dots, x_n 에 기초하여 반파라메터최량판별 함수 $F(x; \theta^*, \mathbf{g}^*)$ 의 미지파라메터 θ^* 과 미지함수 \mathbf{g}^* 을 추정하는 한가지 방법을 논의한다.

먼저 반파라메터최량판별 함수를 추정하자.

표본자료 x_1, \dots, x_n 에 기초하여 반파라메터판별 함수 $F(x; \theta^*, \mathbf{g}^*)$ 의 θ^* 과 \mathbf{g}^* 을 추정하기 위하여 기준우연량 $p_k = \begin{cases} 1, & x_k \in G_1 \\ 0, & x_k \in G_2 \end{cases} \quad (k=1, \dots, n)$ 를 도입하면 반파라메터모형

$$p_k = F(x_k; \theta^*, \mathbf{g}^*) + \varepsilon_k \quad (3)$$

$E\varepsilon_k = 0, \text{Var}\varepsilon_k = P_k \cdot Q_k, \quad P_k = P\{G_1 \mid x_k\} = F(x_k; \theta^*, \mathbf{g}^*), \quad Q_k = 1 - P_k \quad (k=1, \dots, n)$ 가 얻어진다. 여기서 $E p_k = P_k, \text{Var} p_k = (1 - P_k)^2 \cdot P_k + (0 - P_k)^2 \cdot Q_k = P_k \cdot Q_k$ 이다.

이때 $F(x; \theta^*, \mathbf{g}^*)$ 에서 x 를 고정하면 F 는

$$\mathbf{a}^* = (\theta^{*\top}, \mathbf{g}^{*\top})^\top \in \Theta \times H = \mathbf{R}^{s+q}$$

에서 정의된 함수로 볼수 있다.

논문에서는 F 가 \mathbf{a}^* 에 관하여 미분가능한 경우만을 논의한다.

\mathbf{R}^{s+q} 에서 정의된 $F(x; \mathbf{a}^*)$ 을 $\mathbf{a}_0 = (\theta_0^\top, \mathbf{g}_0^\top)^\top$ 에서 테일러전개하면 다음과 같다.

$$F(x; \mathbf{a}^*) = F(x; \mathbf{a}_0) + F'_{\mathbf{a}^*}(x; \mathbf{a}_0)(\mathbf{a}^* - \mathbf{a}_0) + o(\|\mathbf{a}^* - \mathbf{a}_0\|)$$

여기서 $o(\|\mathbf{a}^* - \mathbf{a}_0\|) = r_2(\mathbf{a}_0, l) = r_2(x_k; \mathbf{a}_0, l) \equiv R(x)$ ($R(x)$ 는 θ^*, \mathbf{g}^* 에 관계되는 반파라메터함수이다.)이며 $l = \mathbf{a}^* - \mathbf{a}_0$ 이다.

따라서 반파라메터모형 (3)은 부분선형모형

$$\tilde{p}_k = F'_{\theta^*}(x_k; \mathbf{a}_0) \cdot \theta^* + F'_{\mathbf{g}^*}(x_k; \mathbf{a}_0) \cdot \mathbf{g}^* + R(x_k) + \varepsilon_k \quad (4)$$

$$E\varepsilon_k = 0, \text{Var}\varepsilon_k = P_k \cdot Q_k, \quad P_k = F(x_k; \mathbf{a}^*), \quad Q_k = 1 - P_k \\ \tilde{p}_k = p_k - F(x_k; \mathbf{a}_0) + F'_{\mathbf{a}^*}(x_k; \mathbf{a}_0)\mathbf{a}_0 \quad (k=1, \dots, n)$$

으로 되며 초기추정량 $\hat{\mathbf{a}}_1$ 에 대하여 부분선형모형 (4)는 다음과 같다.

$$\mathbf{Y} = \mathbf{X} \cdot \theta^* + \mathbf{Z} \cdot \mathbf{g}^* + \varepsilon \quad (5)$$

$$E\varepsilon = \mathbf{0}_{n \times 1}, \text{Var}\varepsilon \approx \hat{\text{Var}}\varepsilon = \text{diag}(\hat{P}_1 \cdot \hat{Q}_1, \dots, \hat{P}_n \cdot \hat{Q}_n)_{n \times n}$$

$$\hat{P}_k = F(x_k; \hat{\mathbf{a}}_1), \quad \hat{Q}_k = 1 - \hat{P}_k \quad (k=1, \dots, n)$$

$$\mathbf{X} = (X_1, \dots, X_n)^\top = (F'_{\theta^*}(x_1; \hat{\mathbf{a}}_1)^\top, \dots, F'_{\theta^*}(x_n; \hat{\mathbf{a}}_1)^\top)_{n \times s}^\top$$

$$\mathbf{Y} = (Y_1, \dots, Y_n)^\top = (p_1 - F(x_1; \hat{\mathbf{a}}_1) + F'_{\mathbf{a}^*}(x_1; \hat{\mathbf{a}}_1), \dots, p_n - F(x_n; \hat{\mathbf{a}}_1) + F'_{\mathbf{a}^*}(x_n; \hat{\mathbf{a}}_1))_{n \times 1}^\top$$

$$\mathbf{Z} = (Z_1, \dots, Z_n)^\top = ((F'_{\mathbf{g}^*}(x_1; \hat{\mathbf{a}}_1), l)^\top, \dots, (F'_{\mathbf{g}^*}(x_n; \hat{\mathbf{a}}_1), l)^\top)_{n \times (q+1)}^\top$$

$$\bar{\mathbf{g}}^* = (\mathbf{g}^{*\top}, \mathbf{R})^\top_{(q+1) \times 1}, \quad \boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$$

모형 (5)의 미지파라메터 $\boldsymbol{\theta}^*$ 과 미지함수 $\bar{\mathbf{g}}^*$ 을 추정하기 위하여 측면최소두제곱법과 국부선형화수법을 적용하면 다음의 결론이 나온다.

보조정리 모형 (5)에서 오차들의 분산이 같다고 할 때 반파라메터최량판별규칙의 파라메터성분의 추정량은 $\tilde{\boldsymbol{\theta}}'_n = \left(\sum_{k=1}^n \tilde{X}_k \tilde{X}_k^\top \right)^{-1} \left(\sum_{k=1}^n \tilde{X}_k \tilde{Y}_k \right)_{s \times 1}$ 이며 비파라메터성분의 추정량은

$$\tilde{\bar{\mathbf{g}}}(x) = (\tilde{\mathbf{g}}(x), \tilde{R}(x)) = (\tilde{g}_1(x), \dots, \tilde{g}_q(x), \tilde{R}(x))^\top = (\mathbf{E}_{q+1}, \mathbf{0}_{(q+1) \times (q+1)}) (\mathbf{D}_x^\top \mathbf{W}_x \mathbf{D}_x)^{-1} \mathbf{D}_x^\top \mathbf{W}_x (\mathbf{Y} - \mathbf{X} \tilde{\boldsymbol{\theta}}'_n)$$

이다. 여기서 $\tilde{\mathbf{Y}} = (\tilde{Y}_1, \dots, \tilde{Y}_n)^{-1} = \mathbf{Y} - \mathbf{S}\mathbf{Y}$, $\tilde{\mathbf{X}} = (\tilde{X}_1, \dots, \tilde{X}_n)^\top = \mathbf{X} - \mathbf{S}\mathbf{X}$,

$$\mathbf{S} = \begin{pmatrix} (Z_1^\top, \mathbf{0})(\mathbf{D}_{x_1}^\top \mathbf{W}_{x_1} \mathbf{D}_{x_1})^{-1} \mathbf{D}_{x_1}^\top \mathbf{W}_{x_1} \\ \vdots \\ (Z_n^\top, \mathbf{0})(\mathbf{D}_{x_n}^\top \mathbf{W}_{x_n} \mathbf{D}_{x_n})^{-1} \mathbf{D}_{x_n}^\top \mathbf{W}_{x_n} \end{pmatrix}_{n \times n}, \quad \mathbf{D}_x = \begin{pmatrix} Z_1^\top & (x_1 - x)/h \cdot Z_1^\top \\ \vdots & \vdots \\ Z_n^\top & (x_n - x)/h \cdot Z_n^\top \end{pmatrix}_{n \times 2(q+1)}$$

$$\mathbf{W}_x = \text{diag}(K_h(x_1 - x), \dots, K_h(x_n - x))_{n \times n}, \quad (Z_k^\top, \mathbf{0}) = (Z_{k1}, \dots, Z_{kq}, 0, \dots, 0)_{1 \times 2(q+1)} \quad (k=1, \dots, n)$$

또한 $K_h(\cdot) = K(\cdot/h)/h$ 이며 $K(\cdot)$ 은 핵함수, $h = h_n$ 은 평활화파라메터, \mathbf{E}_{q+1} 은 $q+1$ 차 원단위행렬이다.

정리 1 반파라메터최량판별규칙 (2)의 파라메터성분의 추정량은

$$\hat{\boldsymbol{\theta}}_n = \left(\sum_{k=1}^n \gamma_k \tilde{X}_k \tilde{X}_k^\top \right)^{-1} \left(\sum_{k=1}^n \gamma_k \tilde{X}_k \tilde{Y}_k \right) \quad (6)$$

이며 비파라메터성분의 추정량은

$$\hat{\bar{\mathbf{g}}}_n(x) = (\hat{g}_1(x), \dots, \hat{g}_q(x), \hat{R}(x))^\top = (\mathbf{E}_{q+1}, \mathbf{0}_{(q+1) \times (q+1)}) (\mathbf{D}_x^\top \mathbf{W}_x \mathbf{D}_x)^{-1} \mathbf{D}_x^\top \mathbf{W}_x (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\theta}}_n) \quad (7)$$

이다. 여기서 무게 $\gamma_k = 1/(\text{Var} \varepsilon_k)$ ($k=1, \dots, n$) 이다.

일반적으로 $\text{Var} \varepsilon_k$ ($k=1, \dots, n$) 는 미지이므로 추정하여 얻은 $\hat{\gamma}_k = 1/\hat{\text{Var}} \varepsilon_k$ ($k=1, \dots, n$) 을 리용하면 추정량 (6), (7)은 다음과 같다.

$$\hat{\boldsymbol{\theta}}_n = \left(\sum_{k=1}^n \hat{\gamma}_k \tilde{X}_k \tilde{X}_k^\top \right)^{-1} \left(\sum_{k=1}^n \hat{\gamma}_k \tilde{X}_k \tilde{Y}_k \right)$$

$$\hat{\bar{\mathbf{g}}}_n(x) = (\hat{g}_1(x), \dots, \hat{g}_q(x), \hat{R}(x))^\top = (\mathbf{E}_{q+1}, \mathbf{0}_{(q+1) \times (q+1)}) (\mathbf{D}_x^\top \mathbf{W}_x \mathbf{D}_x)^{-1} \mathbf{D}_x^\top \mathbf{W}_x (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\theta}}_n)$$

증명 오차분산이 상수가 아닌것을 고려하여 무게붙은 최소두제곱법을 적용하면 보조정리에 의하여 식 (6), (7)로 표시되는 추정량을 얻게 된다.(증명끝)

$F(x; \hat{\boldsymbol{\theta}}, \hat{\mathbf{g}})$ 에 의한 판별규칙 (2)에서 오판별확률은 다음과 같이 표시된다.

$$P = q_1^0 \int_{R_2} f_1^0(x) dx + q_2^0 \int_{R_1} f_2^0(x) dx, \quad R_1 = \{x \mid 2F(x; \hat{\boldsymbol{\theta}}, \hat{\mathbf{g}}) > 1\}, \quad R_2 = \{x \mid 2F(x; \hat{\boldsymbol{\theta}}, \hat{\mathbf{g}}) \leq 1\}$$

다음으로 반파라메터최량판별함수추정량의 성질에 대하여 논의하자.

다음과 같은 조건들에 대하여 보자.

① x 의 밀도함수 $f(x)$ 는 2계련속미분가능하고 $f(x) > 0$ 을 만족시킨다.

② $\forall x \in [0, 1], \Gamma(x) = E(Z_1 Z_1^\top \mid x_1 = x)$ 는 불퇴화행렬이고

$$E(X_1 X_1^\top \mid x_1 = x), \quad \Gamma(x), \quad \Phi(x) = E(Z_1 X_1^\top \mid x_1 = x)$$

는 모두 립쉬츠련속이다. 특히 $\Sigma := (E(\gamma_1(X_1 - \Phi(x_1))^T \Gamma(x_1)^{-1} Z_1)^{\otimes 2})^{-1}$ 은 정값행렬이다.

$$\textcircled{3} \exists t > 2, E \|X_1\|^{2t} < \infty, E \|Z_1\|^{2t} < \infty, \exists s < 2 - t^{-1}, n^{2s-1} h_n \rightarrow \infty \quad (n \rightarrow \infty)$$

$$\textcircled{4} g_j(\cdot) \in C^2[0, 1] \quad (j=1, \dots, q+1) \quad nh_n^8 \rightarrow 0, nh_n^8/(\log n)^2 \rightarrow \infty$$

$\textcircled{5}$ 핵함수 $K(\cdot)$ 은 유계받침 $[-10, 1]$ 을 가지는 대칭밀도함수이다.

$$\textcircled{6} \exists M_1, M_2; 0 < M_1 < \sigma_i^2 < M_2 < +\infty, 1 \leq i \leq n \quad (\text{거의})$$

$$\textcircled{7} \max_{1 \leq i \leq n} |\hat{\gamma}_i - \gamma_i| = o_p(n^{-1/4})$$

정리 2 조건 $\textcircled{1}-\textcircled{6}$ 을 만족시키면 식 (6)으로 표시되는 $\hat{\theta}_n$ 은 θ 의 점근정규추정량이다. 즉 $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N(0, \Sigma)$ 가 성립된다.

정리 3 조건 $\textcircled{1}-\textcircled{7}$ 을 만족시키면 식 (7)로 표시되는 $\hat{\mathbf{g}}_n(x) = (\hat{g}_1(x), \dots, \hat{g}_{q+1}(x))^T$ 에 대하여 $\max_{1 \leq j \leq q+1} \sup_{x \in (0, 1)} |\hat{g}_j(x) - g_j(x)| = o_p(c_n)$, $c_n = (\log h^{-1}/nh)^{1/2} + h^2$ 이 성립된다.

얻어진 추정량을 다시 초기추정량 $\hat{\alpha}_1$ 로 놓고 반복한다.

참 고 문 헌

- [1] G. Melnichenko; J. Multi. Anal., 101, 68, 2010.
- [2] M. Mojirsheibani et al.; J. Multi. Anal., 98, 1051, 2007.
- [3] M. S. Srivastava; J. Multi. Anal., 97, 2057, 2006.
- [4] S. Velilla; J. Multi. Anal., 101, 1239, 2010.
- [5] M. M. Seyam et al.; International Journal of Statistics and Probability, 2, 3, 96, 2013.
- [6] D. R. Jeske et al.; J. Multi. Anal., 101, 1622, 2010.
- [7] Jian Qing Fan et al.; Bernoulli, 11, 6, 1031, 2005.
- [8] T. Otsu; J. Multi. Anal., 98, 1923, 2007.
- [9] Jin Hong You et al.; J. Multi. Anal., 101, 1079, 2010.
- [10] B. G. Manjunath et al.; J. Multi. Anal., 103, 107, 2012.
- [11] Qing Mai et al.; J. Multi. Anal., 135, 175, 2015.

주체107(2018)년 3월 10일 원고접수

The Estimation for Semiparametric Optimal Discriminant Functions

Rim Chang Ho

In this paper, by using samples x_1, \dots, x_n we estimated the unknown parameter θ^* and the unknown function \mathbf{g}^* of semiparametric optimal discriminant function $F(x; \theta^*, \mathbf{g}^*)$ and obtained statistical properties of the estimators.

Key words: semiparametric discriminant function, misclassification probability