

Estimable functions of mixed models

Jaesung Choi^{a,1}

^aDepartment of Statistics, Keimyung University

(Received October 27, 2015; Revised December 22, 2015; Accepted January 19, 2016)

Abstract

This paper discusses how to establish estimable functions when there are fixed and random effects in design models. It proves that estimable functions of mixed models are not related to random effects. A fitting constants method is used to obtain sums of squares due to random effects and Hartley's synthesis is used to calculate coefficients of variance components. To test about the fixed effects the degrees of freedom associated with divisor are determined by means of the Satterthwaite approximation.

Keywords: mixed model, estimable function, fitting constants method, Type I sum of squares, synthesis

1. 서론

실험자료의 분석에 이용되는 고정효과모형의 모수들은 일반적으로 추정가능하지 않다. 이는 자료로부터 추정가능한 모수들의 수보다 더 많은 모수들이 모형에 포함되어 있기 때문이다. 최소제곱법으로 모형내 모수를 추정할 때 추정가능한 모수의 수가 제한되어 있음으로 인해 선형적 제약조건하에서 구해지는 모수의 추정값은 유일하게 결정되지 않는다. 고정효과모형에서 모수들의 추정가능한 함수는 중요하게 취급되고 많은 문헌들이 추정가능함수들의 구성과 특성에 관하여 논의하고 있다. 추정가능함수에 관한 구체적 논의는 Searle 등 (1971), Graybill (1976) 그리고 Milliken과 Johnson (1984) 등에서 살펴볼 수 있다.

고정효과모형의 추정가능함수와 관련된 논의로 추정가능한 함수의 구성과 확인, 추정가능한 함수들의 기저집합 그리고 추정가능함수들의 추론 등이 다루어져 왔다. 모형내 모수 또는 모수들의 함수에 대한 추정가능성(estimability)은 모수에 대한 추론이 행해지기 전에 필수적으로 확인되어야 한다. 추정가능하지 않은 모수 또는 모수들의 함수에 관한 추론은 의미가 없다. 고정효과모형의 추정가능함수는 불편추정량이 존재하는 모수들의 선형함수로 정의되거나 모수들의 함수에 대한 추정값이 최소제곱해의 선택에 상관없이 일정한 값으로 주어지는 모수들의 함수로 정의되기도 한다. 즉, 모수들의 함수로 주어지는 불편추정량이 존재하거나 최소제곱해에 불변인 최소제곱추정치를 갖는 함수로 간주된다.

고정효과모형의 행렬표현식에서 모수벡터의 계수행렬은 일반적으로 불완전계수의 모형행렬로 주어진다. 불완전계수의 모형행렬을 갖는 모수벡터의 해를 구하는 최소제곱방법은 벡터공간에서 사영을 구하는 방법과 동일하다. 고정효과모형에서 추정가능함수와 관련된 사영의 활용에 관한 논의는 Choi (2011, 2012, 2014)에서 보여진다.

¹Department of Statistics, Keimyung University, 1095 Dalgubul-Daero, Dalseogu-Gu, Daegu 42601, Korea.
E-mail: jschoi@kmu.ac.kr

실험의 반응에 영향을 주는 요인으로 고정요인외에 확률요인이 추가될 때 자료를 분석하기 위한 모형은 혼합모형으로 주어진다. 요인들의 수와 요인들의 관계가 교차요인인 가 또는 지분요인인 가에 따라 혼합효과모형도 다양한 형태로 주어질 수 있다. 그러나 혼합효과모형은 고정효과들로 주어지는 고정성분과 확률효과로 주어지는 확률성분 그리고 오차항으로 주어지는 세 성분을 포함하게 된다. 혼합모형에서 고정효과들의 선형함수로 주어지는 추정가능함수와 관련된 문제를 다루어 보기로 한다.

본 연구에서 다루고자 하는 부분은 고정효과모형에서 정의되는 추정가능함수에 관한 논의가 혼합효과모형에서도 동일하게 적용될 수 있는 가를 살펴보고 고정효과들의 선형함수가 추정가능함수인 가에 관한 판단이 어떻게 이루어져야 하는 가에 두고 있다. 추정가능함수의 추론과 관련하여 추정가능함수에 관한 추정량의 분산이 분산성분들의 선형함수로 표현되기 때문에 추정량의 분산을 추정하기 위한 방법으로 Satterthwaite의 근사적 방법을 다루게 된다. 그리고 확률효과에 따른 분산성분의 추정을 위한 잔차 모형의 적합에서 사영제곱합을 구하는 방법과 제곱합의 기댓값 계산방법을 다루고 있다.

2. 혼합모형의 추정가능함수

일반적인 혼합모형은 실험단위의 반응벡터를 \mathbf{y} 라 둘 때

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta} + \mathbf{X}_2\boldsymbol{\delta} + \boldsymbol{\epsilon} \quad (2.1)$$

이다. 단, \mathbf{y} 는 크기가 $n \times 1$ 인 관측벡터이고 \mathbf{X}_1 은 $n \times p$ 인 계수행렬이다. \mathbf{X}_1 의 계수는 k ($k < p$)로 가정한다. $\boldsymbol{\beta}$ 는 크기가 $p \times 1$ 인 모수벡터이다. \mathbf{X}_2 는 $n \times r$ 인 계수행렬이고 $\boldsymbol{\delta}$ 는 $r \times 1$ 인 q 개 확률효과와 분산성분과 관련된 확률효과벡터로 크기가 $r \times 1$ 이다. 즉, $\boldsymbol{\delta} = (\boldsymbol{\delta}_1, \boldsymbol{\delta}_2, \dots, \boldsymbol{\delta}_q)'$ 이고 $r = r_1 + r_2 + \dots + r_q$ 이다. 크기가 $r_j \times 1$ 인 $\boldsymbol{\delta}_j$ 는 서로 독립이고 $N(\mathbf{0}, \sigma_j^2 \mathbf{I}_{r_j})$ 인 분포를 가정한다. $\boldsymbol{\epsilon}$ 은 $n \times 1$ 인 오차벡터이고 $N(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I})$ 인 분포를 가정한다. 관측벡터 \mathbf{y} 의 분산은

$$\begin{aligned} \boldsymbol{\Sigma} &= \text{Var}(\mathbf{y}) \\ &= \mathbf{X}_2 \text{Var}(\boldsymbol{\delta}) \mathbf{X}_2' + \sigma_\epsilon^2 \mathbf{I} \end{aligned} \quad (2.2)$$

이다. 단, $\mathbf{X}_2 = (\mathbf{X}_{21}, \mathbf{X}_{22}, \dots, \mathbf{X}_{2q})$ 이다. 식 (2.1)의 혼합모형에서 고정효과들의 추정가능함수를 생각해 보기로 한다. 추정가능함수에 대한 정의가 고정효과모형에서 이루어졌으나 혼합모형에서도 고정효과들의 함수로 주어지는 추정가능한 함수에 대한 추론을 위해 추정가능함수에 관한 유형과 성질을 파악하는 것이 중요하다. 추정가능함수에 대한 Searle 등 (1971)의 정의는 모수들의 한 선형함수가 관측벡터 \mathbf{y} 의 기댓값의 어떤 함수와 일치하면 추정가능하다고 정의한다. 즉, 어떤 벡터 \mathbf{a}' 에 대해 $\mathbf{l}'\boldsymbol{\beta} = \mathbf{a}'E(\mathbf{y})$ 이면 $\mathbf{l}'\boldsymbol{\beta}$ 는 추정가능함을 의미한다. 다시말하면 $\mathbf{l}'\boldsymbol{\beta} = \mathbf{a}'E(\mathbf{y})$ 이 되는 \mathbf{a}' 이 존재하면 $\mathbf{l}'\boldsymbol{\beta}$ 는 추정가능함수로 간주된다. 추정가능함수의 정의는 단순히 \mathbf{y} 의 기댓값의 한 선형함수 $\mathbf{a}'E(\mathbf{y})$ 가 $\mathbf{l}'\boldsymbol{\beta}$ 되는 \mathbf{a}' 이 존재함을 보여주는 것으로 만족된다. $\mathbf{a}'E(\mathbf{y}) = E(\mathbf{a}'\mathbf{y})$ 이므로 $\mathbf{l}'\boldsymbol{\beta}$ 되는 \mathbf{y} 의 많은 선형함수 중 임의 하나로 추정가능성을 입증하는 데 충분하다.

식 (2.1)에서 모수벡터 $\boldsymbol{\beta}$ 의 한 선형함수 $\mathbf{l}'\boldsymbol{\beta}$ 가 추정가능함수이기 위해서는 $E(\mathbf{a}'\mathbf{y})$ 가 $\mathbf{l}'\boldsymbol{\beta}$ 되는 \mathbf{a}' 가 존재함을 보여야 한다. $\boldsymbol{\beta}$ 는 모형행렬 \mathbf{X}_1 의 행공간에서 한 좌표를 나타내는 모수점이다. 따라서 $\mathbf{X}_1^- \mathbf{X}_1 \boldsymbol{\beta} = \boldsymbol{\beta}$ 인 성질을 만족한다. 여기서 \mathbf{X}_1^- 는 Moore-Penrose의 일반화된 역행렬(generalized inverse)을 나타낸다. $\mathbf{l}'\boldsymbol{\beta}$ 가 모수벡터 $\boldsymbol{\beta}$ 의 선형함수로 추정가능하면

$$\mathbf{l}' \mathbf{X}_1^- \mathbf{X}_1 \boldsymbol{\beta} = \mathbf{l}'\boldsymbol{\beta} \quad (2.3)$$

가 만족된다. 따라서

$$\mathbf{l}' \mathbf{X}_1^- \mathbf{X}_1 = \mathbf{l}' \quad (2.4)$$

이다. 식 (2.4)는 모수벡터 β 의 한 선형함수가 추정가능한 함수인가의 판단을 위한 계수벡터 l' 의 조건을 나타내고 있다. $l'\beta$ 되는 a' 은

$$\begin{aligned} a'E(y) &= a'X_1\beta & (2.5) \\ &= (a'X_1)X_1^-X_1\beta \\ &= l'\beta \end{aligned}$$

이므로 $a'X_1 = l'$ 임을 알 수 있다. X_1 이 비정칙이므로 l' 되는 많은 해가 존재하게 되고 그 중 임의의 a' 를 구할 수 있다. 고정효과모형에서 정의된 추정가능함수 $l'\beta$ 는 불편추정량이 존재하는 β 의 선형함수이어야 하므로 l' 은 X_1 의 행공간에 존재하는 벡터가 되어야 함을 보였다. 이는 모수벡터 β 의 추정가능함수를 얻기 위한 계수벡터 l' 는 X_1 의 행벡터들의 한 선형결합으로 얻을 수 있음을 의미하고 있다. 크기가 $n \times p$ 인 모형행렬 X_1 의 계수가 k 이므로 k 개의 선형독립인 추정가능함수들을 구할 수 있고 이들의 집합이 β 의 한 기저집합으로 정의된다. 추정가능함수들의 한 기저집합은 X_1X_1' 의 기약가우스 행렬을 이용하는 등의 다양한 방법으로 구할 수 있다. 이러한 논의는 Choi (2012)와 Elswick 등 (1991)에서 볼 수 있다.

모수벡터 β 의 추정가능함수의 형태는 계획행렬 X_1 의 행벡터들의 선형결합에 의해 결정되므로 β 의 추정량이 불편추정량이기만 하면 추정가능함수의 정의를 만족하게 된다. 불편성의 성질을 만족시키는 고정효과의 선형함수는 혼합모형에서도 추정가능함을 살펴보기로 한다.

3. 분산성분의 추정

혼합모형은 고정효과의 고정성분과 확률효과와 오차로 주어지는 확률성분으로 구분될 수 있다. 혼합모형의 분석은 고정효과의 추론에 앞서 분산성분의 추론을 위한 확률모형으로 시작한다. 확률모형은 모형의 고정효과부분을 적합시켜 얻은 잔차에 대한 모형으로 주어진다. 고정효과가 제외된 확률효과모형은 식 (2.1)의 고정성분 $X_1\beta$ 를 적합시켜 구해진 잔차벡터를 r 이라 둘 때 $r = y - X_1\hat{\beta}$ 으로 주어진다. Henderson (1953)의 방법 3 또는 상수적합법에 의한 고정효과의 적합에 따른 제곱합은 확률성분의 분산성분과 무관하게 구해진다. 분산성분의 추정에 관한 논의는 Corbeil과 Searle (1976) 그리고 Searle 등 (1992)에서 보여진다. r 의 모형은

$$\begin{aligned} r &= (I - X_1X_1^-)y & (3.1) \\ &= (I - X_1X_1^-)X_2\delta + \epsilon_r \end{aligned}$$

이다. 단, ϵ_r 은 $(I - X_1X_1^-)\epsilon$ 이다. 확률벡터 δ 의 분산성분벡터를 σ_δ^2 라 두면 $\sigma_\delta^2 = (\sigma_1^2, \sigma_2^2, \dots, \sigma_q^2)'$ 이다. 확률벡터 δ 의 분산성분 σ_j^2 ($j = 1, 2, \dots, q$)을 추정하기 위해 모형의 단계별 적합에 의한 확률모형을 이용하게 된다. 단계별 적합의 확률모형에서 성분벡터 δ_j 에 따른 모형행렬을 X_{2j} 라 두면 X_{2j} 로의 사영을 이용한 제곱합과 제곱합의 기댓값을 Hartley (1967)의 합성법(synthesis)으로 구할 수 있다. 식 (3.1)로부터 확률벡터 δ_1 에 따른 제곱합을 구하기 위한 잔차모형은

$$\begin{aligned} r &= (I - X_1X_1^-)X_{21}\delta_1 + (I - X_1X_1^-)(X_{22}\delta_2 + \dots + X_{2q}\delta_q + \epsilon_r) & (3.2) \\ &= (I - X_1X_1^-)X_{21}\delta_1 + \epsilon_1 \end{aligned}$$

으로 주어진다. 단, $\epsilon_1 = (I - X_1X_1^-)(X_{22}\delta_2 + \dots + X_{2q}\delta_q + \epsilon_r)$ 이다. 식 (3.2)의 모형행렬을 M_1 이라 두자. $M_1 = (I - X_1X_1^-)X_{21}$ 이고 M_1 으로의 사영은 $M_1M_1^-r$ 이고 사영까지의 거리제곱합은

$\mathbf{r}'\mathbf{M}_1\mathbf{M}_1^-\mathbf{r}$ 이다. δ_1 에 따른 제곱합을 Q_1 이라 두면 $Q_1 = \mathbf{r}'\mathbf{M}_1\mathbf{M}_1^-\mathbf{r}$ 로 구해진다. 식 (3.2)로부터 확률벡터 δ_2 에 따른 제곱합을 구하기 위한 잔차모형을 \mathbf{r}_1 이라 두면

$$\begin{aligned}\mathbf{r}_1 &= (\mathbf{I} - \mathbf{X}_1\mathbf{X}_1^- - \mathbf{M}_1\mathbf{M}_1^-)\mathbf{X}_{22}\delta_2 + (\mathbf{I} - \mathbf{X}_1\mathbf{X}_1^- - \mathbf{M}_1\mathbf{M}_1^-)(\mathbf{X}_{23}\delta_3 + \cdots + \mathbf{X}_{2q}\delta_q + \epsilon_1) \\ &= (\mathbf{I} - \mathbf{X}_1\mathbf{X}_1^- - \mathbf{M}_1\mathbf{M}_1^-)\mathbf{X}_{22}\delta_2 + \epsilon_2\end{aligned}\quad (3.3)$$

으로 주어진다. 식 (3.3)의 모형행렬을 \mathbf{M}_2 라 두자. $\mathbf{M}_2 = (\mathbf{I} - \mathbf{X}_1\mathbf{X}_1^- - \mathbf{M}_1\mathbf{M}_1^-)\mathbf{X}_{22}$ 이고 \mathbf{M}_2 로의 사영은 $\mathbf{M}_2\mathbf{M}_2^-\mathbf{r}_1$ 이고 사영까지의 거리제곱합은 $\mathbf{r}_1'\mathbf{M}_1\mathbf{M}_1^-\mathbf{r}_1$ 이다. δ_2 에 따른 제곱합을 Q_2 이라 두면 $Q_2 = \mathbf{r}_1'\mathbf{M}_2\mathbf{M}_2^-\mathbf{r}_1$ 로 구해진다. 분산성분 σ_q^2 과 관련된 확률벡터 δ_q 에 따른 제곱합을 구하기 위한 잔차모형을 \mathbf{r}_q 라 두면

$$\mathbf{r}_{q-1} = (\mathbf{I} - \mathbf{X}_1\mathbf{X}_1^- - \mathbf{M}_1\mathbf{M}_1^- - \cdots - \mathbf{M}_{q-1}\mathbf{M}_{q-1}^-)\mathbf{X}_{2q}\delta_q + \epsilon_q \quad (3.4)$$

으로 주어진다. 단, $\epsilon_q = (\mathbf{I} - \mathbf{X}_1\mathbf{X}_1^- - \mathbf{M}_1\mathbf{M}_1^- - \cdots - \mathbf{M}_{q-1}\mathbf{M}_{q-1}^-)\epsilon_{q-1}$ 이다. 식 (3.4)의 모형행렬을 \mathbf{M}_q 라 두자. $\mathbf{M}_q = (\mathbf{I} - \mathbf{X}_1\mathbf{X}_1^- - \mathbf{M}_1\mathbf{M}_1^- - \cdots - \mathbf{M}_{q-1}\mathbf{M}_{q-1}^-)\mathbf{X}_{2q}$ 이고 \mathbf{M}_q 로의 사영은 $\mathbf{M}_q\mathbf{M}_q^-\mathbf{r}_{q-1}$ 이고 사영까지의 거리제곱합은 $\mathbf{r}_{q-1}'\mathbf{M}_1\mathbf{M}_1^-\mathbf{r}_{q-1}$ 이다. δ_q 에 따른 제곱합을 Q_q 이라 두면 $Q_q = \mathbf{r}_{q-1}'\mathbf{M}_q\mathbf{M}_q^-\mathbf{r}_{q-1}$ 로 구해진다.

확률요인에 따른 변동량들은 Q_1, Q_2, \dots, Q_q 로 구해지는 q 개의 제곱합이고 이들은 반응벡터 \mathbf{y} 의 이차형식으로 표현된다. 분산성분의 추정값들은 이들 변동량과 기댓값을 같게 두는 연립방정식의 해로 주어지기 때문에 분산성분들의 적률추정치가 구해지면 고정효과의 선형함수이면서 추정가능함수인 $\mathbf{l}'\beta$ 의 불편추정량을 구할 수 있다. $\mathbf{l}' = \mathbf{l}'\mathbf{X}_1^-\mathbf{X}_1$ 되는 $\mathbf{a}'\mathbf{X}_1$ 의 \mathbf{a}' 은 \mathbf{X}_1 의 행벡터들의 선형결합으로 얻을 수 있는 임의의 한 계수벡터이므로 β 의 선형함수 $\mathbf{l}'\beta$ 의 추정가능성을 보여주는 데 이용된다. 추정가능함수 $\mathbf{l}'\beta$ 의 불편추정량을 구하는 문제를 생각해 보기로 한다.

4. 추정가능함수의 추정

혼합모형의 고정효과벡터에 대한 추정방법으로 가중최소제곱추정량을 이용해 보기로 한다. 식 (2.1)의 모형에 상수적합법을 이용하여 구한 식 (2.2)의 공분산행렬 Σ 의 추정행렬을 $\hat{\Sigma}$ 라 두자. 추정된 공분산행렬 $\hat{\Sigma}$ 을 이용한 모수벡터 β 의 추정량을 $\hat{\beta}$ 이라 두자. $\hat{\beta}$ 은 가중최소제곱법에 의해

$$\hat{\beta} = (\mathbf{X}_1'\hat{\Sigma}^{-1}\mathbf{X}_1)^- \mathbf{X}_1'\hat{\Sigma}^{-1}\mathbf{y} \quad (4.1)$$

로 구해진다. $\mathbf{l}' = \mathbf{l}'\mathbf{X}_1^-\mathbf{X}_1$ 되는 $\mathbf{l}'\beta$ 의 추정량을 $\mathbf{l}'\hat{\beta}$ 라 두면

$$\begin{aligned}E(\mathbf{l}'\hat{\beta}) &= \mathbf{l}'E(\hat{\beta}) \\ &= \mathbf{l}'(\mathbf{X}_1'\hat{\Sigma}^{-1}\mathbf{X}_1)^- \mathbf{X}_1'\hat{\Sigma}^{-1}E(\mathbf{y}) \\ &= \mathbf{l}'\beta\end{aligned}\quad (4.2)$$

이다. $\mathbf{l}'\hat{\beta}$ 이 $\mathbf{l}'\beta$ 의 불편추정량이고 β 의 선형함수이므로 추정가능함수임을 나타내고 있다. $\mathbf{l}'\beta$ 가 추정가능함수일 때 불편추정량을 $\mathbf{l}'\hat{\beta}$ 라 두자. 식 (2.1)의 모형에 대한 가정으로부터 $\mathbf{l}'\hat{\beta}$ 은 평균이 $\mathbf{l}'\beta$ 이고 분산은

$$\begin{aligned}\text{Var}(\mathbf{l}'\hat{\beta}) &= \mathbf{l}'(\mathbf{X}_1'\Sigma^{-1}\mathbf{X}_1)^- \mathbf{l} \\ &= \mathbf{l}'(\mathbf{X}_1'[\mathbf{X}_2\text{Var}(\delta)\mathbf{X}_2' + \sigma_\epsilon^2\mathbf{I}]^{-1}\mathbf{X}_1)^- \mathbf{l}\end{aligned}\quad (4.3)$$

$$= \mathbf{l}' \left[\mathbf{X}'_1 (\sigma_1^2 \mathbf{X}_{21} \mathbf{X}'_{21} + \sigma_2^2 \mathbf{X}_{22} \mathbf{X}'_{22} + \cdots + \sigma_q^2 \mathbf{X}_{2q} \mathbf{X}'_{2q} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{X}_1 \right]^{-1} \mathbf{l}$$

이다. 추정가능함수 $\mathbf{l}'\boldsymbol{\beta}$ 의 추정량 $\mathbf{l}'\hat{\boldsymbol{\beta}}$ 의 분포는 $N(\mathbf{l}'\boldsymbol{\beta}, \text{Var}(\mathbf{l}'\hat{\boldsymbol{\beta}}))$ 이다. $\text{Var}(\mathbf{l}'\hat{\boldsymbol{\beta}})$ 가 미지이므로 분산의 추정값을 $\hat{\sigma}_i^2$ 으로 나타내면 식 (4.3)으로부터

$$\hat{\sigma}_i^2 = \mathbf{l}' \left[\mathbf{X}'_1 (\hat{\sigma}_1^2 \mathbf{X}_{21} \mathbf{X}'_{21} + \hat{\sigma}_2^2 \mathbf{X}_{22} \mathbf{X}'_{22} + \cdots + \hat{\sigma}_q^2 \mathbf{X}_{2q} \mathbf{X}'_{2q} + \hat{\sigma}_\epsilon^2 \mathbf{I})^{-1} \mathbf{X}_1 \right]^{-1} \mathbf{l} \quad (4.4)$$

로 추정된다. 식 (4.4)에서 분산성분 σ_i^2 의 추정값 $\hat{\sigma}_i^2$ 은 \mathbf{y} 의 이차형식들의 선형결합으로 구해진다. 이들 이차형식은 변동요인에 따른 자유도를 갖기 때문에 $\hat{\sigma}_i^2$ 의 자유도를 Satterthwaite (1946)에 의해 근사적으로 구할 수 있다. $\hat{\sigma}_i^2$ 의 자유도를 구하기 위한 추정가능함수 $\mathbf{l}'\hat{\boldsymbol{\beta}}$ 의 분산식 (4.3)은

$$\begin{aligned} \text{Var}(\mathbf{l}'\hat{\boldsymbol{\beta}}) &= g(\sigma_1^2, \sigma_2^2, \dots, \sigma_q^2, \sigma_\epsilon^2) c_l \\ &= E(MQ_l) c_l \end{aligned} \quad (4.5)$$

의 형태로 표현된다. 단, $g(\sigma_1^2, \sigma_2^2, \dots, \sigma_q^2, \sigma_\epsilon^2)$ 는 분산성분들의 선형결합으로 주어지는 함수이다. c_l 은 상수이고 MQ_l 은 $E(MQ_l) = g(\sigma_1^2, \sigma_2^2, \dots, \sigma_q^2, \sigma_\epsilon^2)$ 이되는 추정량이고 $MQ_l = k_1 MQ_1 + k_2 MQ_2 + \cdots + k_q MQ_q + k_\epsilon MQ_\epsilon$ 로 구해진다. MQ_i 는 자유도가 df_i 인 Q_i 의 평균제곱을 나타낸다. MQ_i 의 자유도 df_i 은 Satterthwaite의 근사적 자유도를 이용하여 구해진다. MQ_i 의 자유도 df_i 은

$$df_l = \frac{(MQ_l)^2}{(k_1 MQ_1)^2/df_1 + (k_2 MQ_2)^2/df_2 + \cdots + (k_q MQ_q)^2/df_q + (k_\epsilon MQ_\epsilon)^2/df_\epsilon} \quad (4.6)$$

로 구해진다. 추정가능함수 $\mathbf{l}'\boldsymbol{\beta}$ 의 $100(1-\alpha)\%$ 신뢰구간은 $(\mathbf{l}'\hat{\boldsymbol{\beta}} - t_{\alpha/2, df_l} \sqrt{\hat{\sigma}_i^2}, \mathbf{l}'\hat{\boldsymbol{\beta}} + t_{\alpha/2, df_l} \sqrt{\hat{\sigma}_i^2})$ 이다.

5. 건조막 자료의 예

Table 5.1은 Hicks (1973)의 건조막 두께(dry-film thickness)에 대한 실험자료이다. 건조막 두께 측정(Y)에 영향을 주는 요인으로 측정일, 측정기사 그리고 문설정의 세 요인을 생각한다. 실험에서 측정일(D)는 실험이 예정된 시기의 그 달에 임의로 2일이 선정되고 세 고정수준의 문설정(G)에 광택액의 건조막 측정을 위한 측정기사(O)는 측정기사들의 모집단에서 임의로 세 사람이 선정되어 행해진 실험으로 요인 D 와 O 는 확률요인이고 G 는 고정요인이다. D 는 1, 2의 두 수준이고 O 는 o_1, o_2, o_3 의 세 수준이며 G 는 g_1, g_2, g_3 의 세 수준이다. 세 요인의 결합수준에서 2회 측정된 자료를 나타낸다. 건조막 두께의 측정자료를 분석하기 위한 모형은

$$\begin{aligned} \mathbf{y} &= \mathbf{j}\mu + \mathbf{X}_d \boldsymbol{\alpha} + \mathbf{X}_o \boldsymbol{\gamma} + \mathbf{X}_{do}(\boldsymbol{\alpha}\boldsymbol{\gamma}) + \mathbf{X}_g \boldsymbol{\tau} + \mathbf{X}_{dg}(\boldsymbol{\alpha}\boldsymbol{\tau}) + \mathbf{X}_{og}(\boldsymbol{\gamma}\boldsymbol{\tau}) + \mathbf{X}_{dog}(\boldsymbol{\alpha}\boldsymbol{\gamma}\boldsymbol{\tau}) + \boldsymbol{\epsilon} \\ &= (\mathbf{j}\mu + \mathbf{X}_g \boldsymbol{\tau}) + [\mathbf{X}_d \boldsymbol{\alpha} + \mathbf{X}_o \boldsymbol{\gamma} + \mathbf{X}_{do}(\boldsymbol{\alpha}\boldsymbol{\gamma}) + \mathbf{X}_{dg}(\boldsymbol{\alpha}\boldsymbol{\tau}) + \mathbf{X}_{og}(\boldsymbol{\gamma}\boldsymbol{\tau}) + \mathbf{X}_{dog}(\boldsymbol{\alpha}\boldsymbol{\gamma}\boldsymbol{\tau})] + \boldsymbol{\epsilon} \\ &= \mathbf{X}_1 \boldsymbol{\beta} + \mathbf{X}_2 \boldsymbol{\delta} + \boldsymbol{\epsilon} \end{aligned} \quad (5.1)$$

이다. 즉, 식 (2.1)의 세 성분으로 표현됨을 알 수 있다. 혼합효과모형의 분석은 확률효과모형으로 시작하게 된다. 고정효과벡터에 종속되지 않는 확률효과모형은 고정효과를 적합시킨 잔차벡터에 대한 모형이다. 식 (3.1)의 모형으로부터 $\boldsymbol{\alpha}$ 에 따른 제곱합을 얻기 위한 모형은

$$\mathbf{r} = (\mathbf{I} - \mathbf{X}_1 \mathbf{X}_1^{-}) \mathbf{X}_d \boldsymbol{\alpha} + \boldsymbol{\epsilon}_1 \quad (5.2)$$

이다. 식 (5.2)의 적합에서 $\boldsymbol{\alpha}$ 에 따른 변동량을 나타내는 제곱합은 0으로 구해진다. 이때 제곱합을 구하기 위한 사영공간으로의 사영행렬을 $\mathbf{M}_d \mathbf{M}_d^{-}$ 라 두면 해당하는 사영제곱합은 $\mathbf{r}' \mathbf{M}_d \mathbf{M}_d^{-} \mathbf{r}$ 가 된다.

Table 5.1. Dry-film thickness data

Day	Operator	Gate setting	Measurements	Day	Operator	Gate setting	Measurements
1	o_1	g_1	0.38, 0.40	2	o_1	g_1	0.40, 0.40
1	o_2	g_1	0.39, 0.41	2	o_2	g_1	0.39, 0.43
1	o_3	g_1	0.45, 0.40	2	o_3	g_1	0.41, 0.40
1	o_1	g_2	0.63, 0.59	2	o_1	g_2	0.68, 0.66
1	o_2	g_2	0.72, 0.70	2	o_1	g_2	0.77, 0.76
1	o_3	g_2	0.78, 0.79	2	o_1	g_2	0.85, 0.84
1	o_1	g_3	0.76, 0.78	2	o_1	g_3	0.86, 0.82
1	o_2	g_3	0.95, 0.96	2	o_1	g_3	0.86, 0.85
1	o_3	g_3	1.03, 1.06	2	o_1	g_3	1.01, 0.98

$M_d = (I - X_1 X_1^-) X_d$ 로 주어진다. γ 에 따른 변동량을 구하기 위한 잔차벡터를 r_1 이라 둘 때 r_1 에 대한 모형은

$$r_1 = (I - X_1 X_1^- - M_d M_d^-) X_o \gamma + \epsilon_2 \quad (5.3)$$

이다. 식 (5.3)에서 $(I - X_1 X_1^- - M_d M_d^-) X_o$ 를 M_o 라 두면 γ 에 따른 제곱합은 $r_1' M_o M_o^- r_1$ 에 의해 구해지며 그 값은 0.1121이다. 요인 D 와 O 의 교호작용 DO 에 따른 변동량을 구하기 위한 잔차벡터를 r_2 라 둘 때 r_2 에 대한 모형은

$$r_2 = (I - X_1 X_1^- - M_d M_d^- - M_o M_o^-) X_{do}(\alpha\gamma) + \epsilon_3 \quad (5.4)$$

이다. 식 (5.4)에서 $(I - X_1 X_1^- - M_d M_d^- - M_o M_o^-) X_{do}$ 를 M_{do} 라 두면 $(\alpha\gamma)$ 에 따른 제곱합은 $r_2' M_{do} M_{do}^- r_2$ 로 구해지고 그 값은 0.0060이다. $(\alpha\tau)$ 에 따른 변동량을 구하기 위한 잔차벡터를 r_3 라 둘 때 r_3 에 대한 모형은

$$r_3 = (I - X_1 X_1^- - M_d M_d^- - M_o M_o^- - M_{do} M_{do}^-) X_{dg}(\alpha\tau) + \epsilon_4 \quad (5.5)$$

이다. 식 (5.5)에서 $(I - X_1 X_1^- - M_d M_d^- - M_o M_o^- - M_{do} M_{do}^-) X_{dg}$ 를 M_{dg} 라 두면 $\alpha\tau$ 에 따른 변동량은 $r_3' M_{dg} M_{dg}^- r_3$ 로 구해지고 그 값은 0.0108로 주어진다. $(\gamma\tau)$ 에 따른 변동량을 구하기 위한 잔차벡터를 r_4 라 둘 때 r_4 에 대한 모형은

$$r_4 = (I - X_1 X_1^- - M_d M_d^- - M_o M_o^- - M_{do} M_{do}^- - M_{dg} M_{dg}^-) X_{og}(\gamma\tau) + \epsilon_5 \quad (5.6)$$

이다. 식 (5.6)에서 $(\gamma\tau)$ 의 계수행렬을 M_{og} 라 두면 $\gamma\tau$ 에 따른 변동량은 $r_4' M_{og} M_{og}^- r_4$ 로 구해지고 그 값은 0.0455로 주어진다. 세 요인의 교호작용 $(\alpha\gamma\tau)$ 에 따른 변동량을 구하기 위한 잔차벡터를 r_5 라 둘 때 r_5 에 대한 모형은

$$r_5 = (I - X_1 X_1^- - M_d M_d^- - M_o M_o^- - M_{do} M_{do}^- - M_{dg} M_{dg}^- - M_{og} M_{og}^-) \times X_{dog}(\alpha\gamma\tau) + \epsilon_6 \quad (5.7)$$

이다. $X_{dog}(\alpha\gamma\tau)$ 의 계수행렬을 M_{dog} 라 두면 $\alpha\gamma\tau$ 에 따른 변동량은 $r_5' M_{dog} M_{dog}^- r_5$ 로 구해지고 그 값은 0.0073으로 주어진다. 오차에 따른 제곱합은 0.0059로 구해진다. 적률법에 의해 7개 분산성분을 추정하기 위한 제곱합을 각기 $Q_d, Q_o, Q_{do}, Q_{dg}, Q_{og}, Q_{dog}$ 와 Q_ϵ 로 나타내면 $Q = (Q_d, Q_o, Q_{do}, Q_{dg}, Q_{og}, Q_{dog}, Q_\epsilon)'$ 인 벡터를 나타낸다.

즉, $Q = (0.0010, 0.1121, 0.0060, 0.0113, 0.0428, 0.0099, 0.0059)'$ 이다. 합성법으로 Q_i 의 기댓값을 구하여 얻은 연립방정식은 교호작용에 제약이 없는 경우에

$$\begin{aligned} Q_d &= 18\sigma_d^2 + 6\sigma_{do}^2 + 6\sigma_{dg}^2 + 2\sigma_{dog}^2 + \sigma_\epsilon^2, \\ Q_o &= 24\sigma_o^2 + 12\sigma_{do}^2 + 8\sigma_{og}^2 + 4\sigma_{dog}^2 + 2\sigma_\epsilon^2, \\ Q_{do} &= 12\sigma_{do}^2 + 4\sigma_{dog}^2 + 2\sigma_\epsilon^2, \\ Q_{dg} &= 12\sigma_{dg}^2 + 4\sigma_{dog}^2 + 2\sigma_\epsilon^2, \\ Q_{og} &= 16\sigma_{og}^2 + 8\sigma_{dog}^2 + 4\sigma_\epsilon^2, \\ Q_{dog} &= 8\sigma_{dog}^2 + 4\sigma_\epsilon^2, \\ Q_\epsilon &= 18\sigma_\epsilon^2 \end{aligned} \tag{5.8}$$

이다. 분산성분들에 대한 적률추정값들은 각기

$$\begin{aligned} \hat{\sigma}_d^2 &= -0.0003, & \hat{\sigma}_o^2 &= 0.0037, & \hat{\sigma}_{do}^2 &= 0.0001, \\ \hat{\sigma}_{dg}^2 &= 0.0005, & \hat{\sigma}_{og}^2 &= 0.0021, & \hat{\sigma}_{dog}^2 &= 0.0011, \\ \hat{\sigma}_\epsilon^2 &= 0.0003 \end{aligned} \tag{5.9}$$

으로 구해진다. 식 (4.1)에 의한 모수벡터 β 의 추정벡터를 구하기 위해 $\hat{\Sigma}$ 은

$$\hat{\Sigma} = \hat{\sigma}_d^2 \mathbf{X}_d \mathbf{X}'_d + \hat{\sigma}_o^2 \mathbf{X}_o \mathbf{X}'_o + \hat{\sigma}_{do}^2 \mathbf{X}_{do} \mathbf{X}'_{do} + \cdots + \hat{\sigma}_{dog}^2 \mathbf{X}_{dog} \mathbf{X}'_{dog} + \hat{\sigma}_\epsilon^2 \mathbf{I} \tag{5.10}$$

이다. $\hat{\sigma}_d^2 = -0.0003$ 이므로 0으로 간주하면

$$\hat{\Sigma} = \hat{\sigma}_o^2 \mathbf{X}_o \mathbf{X}'_o + \hat{\sigma}_{do}^2 \mathbf{X}_{do} \mathbf{X}'_{do} + \cdots + \hat{\sigma}_{dog}^2 \mathbf{X}_{dog} \mathbf{X}'_{dog} + \hat{\sigma}_\epsilon^2 \mathbf{I} \tag{5.11}$$

이다. 식 (4.1)의 $\hat{\beta}$ 은

$$\mathbf{X}'_1 \hat{\Sigma}^{-1} \mathbf{X}_1 \hat{\beta} = \mathbf{X}'_1 \hat{\Sigma}^{-1} \mathbf{y} \tag{5.12}$$

의 가능한 한 해이다. Table 5.1의 자료로부터 구해진 행렬방정식은

$$\begin{pmatrix} 607.8229 & 202.6076 & 202.6076 & 202.6076 \\ 202.6076 & 642.9605 & -220.1764 & -220.1764 \\ 202.6076 & -220.1764 & 642.9605 & -220.1764 \\ 202.6076 & -220.1764 & -220.1764 & 642.9605 \end{pmatrix} \hat{\beta} = \begin{pmatrix} 414.5015 \\ -100.8738 \\ 180.3650 \\ 335.0103 \end{pmatrix} \tag{5.13}$$

이고 $\hat{\beta} = (0.5115, -0.1065, 0.2194, 0.3985)'$ 으로 구해진다. 고정효과벡터 β 에 따른 변동량을 Q_g 라 두면 $Q_g = 1.5732$ 로 구해지고 $E(Q_g) = 12\sigma_{dg}^2 + 8\sigma_{og}^2 + 4\sigma_{dog}^2 + 2\sigma_\epsilon^2 + 24\phi_g$ 로 주어진다. $H_0 : \alpha_1 = \alpha_2 = \alpha_3$ 가 참이면 $\phi_g = 0$ 이므로 MQ_g 를 Q_g 의 평균평방향이라 두자. 검정통계량 $F = MQ_g/MQ_u$ 가 되는 MQ_u 의 기댓값은 $E(MQ_u) = 6\sigma_{dg}^2 + 4\sigma_{og}^2 + 2\sigma_{dog}^2 + \sigma_\epsilon^2$ 이고 $MQ_u = c_1MS_d + c_2MS_o + \cdots + c_\epsilon MS_\epsilon$ 이다. Satterthwaite의 근사화에 의한 자유도는 식 (4.6)에 의해 $r = 5.33$ 으로 계산된다. 자유도 $df_g = 2$ 와 $df_u = 5.33$ 에서 유의수준 0.001의 F -값은 32.93으로 주어지고 계산값이 $F = 80.3778$ 이므로 H_0 를 기각한다. 즉, 문설정(G)의 세 수준 간의 효과는 같지 않음을 보여준다. 모수벡터 β 의 한 선형함수를 $l_2\beta$ 라 두고 $l_2\beta = \mu + \beta_2$ 라 하자. $\mu + \beta_2$ 가 추정가능함수이면 $l_2 = (1, 0, 1, 0)'$ 이 식 (2.4)의 조건이 성립되는 계수벡터이어야 한다. $l'_2 = l_2 \mathbf{X}'_1 \mathbf{X}_1 = (1, 0, 1, 0)$ 이므로 $\mu + \beta_2$ 는 추정가능함수이다.

추정가능함수의 추정값은 식 (5.12)의 한 해인 $\hat{\beta}$ 을 이용하여 $\hat{\mu} + \hat{\beta}_2 = 0.7308$ 을 구할 수 있다. 모수 벡터 β 의 다른 선형함수를 $l_1'\beta$ 라 두고 $l_1'\beta = \mu + \beta_1$ 이라 하자. $l_1' = l_1'X_1^{-1}X_1 = (1, 1, 0, 0)$ 이므로 $\mu + \beta_1$ 은 추정가능함수이다. $\mu + \beta_1$ 의 추정값은 $\hat{\mu} + \hat{\beta}_1 = 0.4050$ 으로 구해진다. 두 추정가능함수 $l_1'\beta = \mu + \beta_1$ 와 $l_2'\beta = \mu + \beta_2$ 의 차이로 주어지는 모수벡터 β 의 함수도 추정가능하다. 즉, $l_1'\beta - l_2'\beta = (l_1' - l_2')\beta$ 이므로 l' 을 $l' = (l_1' - l_2')$ 으로 두면 $l'\beta = \beta_1 - \beta_2$ 이다. 추정가능함수 $l'\beta$ 의 추정치는 $l'\hat{\beta} = -0.3258$ 이고 $l'\hat{\beta}$ 의 분산추정치는 식 (4.4)에 의해 $\hat{\sigma}_l^2 = 0.0023$ 으로 계산된다. $l'\beta$ 의 95% 신뢰 구간은 $(l'\hat{\beta} - t_{0.025, 5.33}\sqrt{\hat{\sigma}_l^2}, l'\hat{\beta} + t_{0.025, 5.33}\sqrt{\hat{\sigma}_l^2})$ 으로 구해진다.

6. 결론

본 논문은 혼합모형에서 고정효과의 한 선형함수가 추정가능함수이기 위한 조건으로 고정효과벡터의 선형결합을 나타내는 계수벡터가 단순히 고정효과벡터의 계수행렬에 의해 생성되는 행공간에 속하면 충분함을 논의하고 있다. 그 근거로 모수벡터 β 와의 선형결합을 나타내는 계수벡터 l' 은 분산성분을 나타내는 확률벡터의 계수행렬과는 아무런 상관이 없음을 구체적으로 증명하고 있다. 이는 고정효과모형에서 정의된 추정가능함수가 혼합모형에서도 적용될 수 있는 가에 대한 의문을 해소시켜 주는 부분으로 생각된다. 혼합효과모형에서 한 선형함수가 추정가능할 때 그 함수의 추론을 위한 방법을 다루고 있으며 고정요인과 확률요인의 효과를 추정하기 위해 상수적합법을 적용하고 있다. 상수적합법은 혼합모형에서 변동요인에 따른 제곱합을 계산하기 위해 이용되고 벡터공간에서의 사영을 이용할 수 있는 이점이 있다. 분산성분과 관련된 각 사영공간으로의 사영행렬이 어떻게 주어지는 가를 모형의 단계별 적합을 이용하여 논의하고 있다. 분산성분을 추정하기 위한 모형의 단계별 적합에서 구해지는 제곱합들은 상호직교하는 벡터부분공간으로의 사영에 의한 제 1종 제곱합임을 보여주며 제 1종 제곱합의 기댓값으로 Hartley의 합성법을 이용하여 계산하고 있다. 또한, 추론을 위한 자유도의 계산에 Satterthwaite의 근사과정을 논의하고 있다.

References

- Choi, J. S. (2011). Type I analysis by projections, *The Korean Journal of Applied Statistics*, **24**, 373–381.
- Choi, J. S. (2012). Type II analysis by projections, *Journal of the Korean Data & Information Science Society*, **23**, 1155–1163.
- Choi, J. S. (2014). Projection analysis for two-way variance components, *Journal of the Korean Data & Information Science Society*, **23**, 547–554.
- Corbeil, R. R. and Searle, S. R. (1976). A comparison of variance component estimators, *Biometrics*, **32**, 779–791.
- Elswick, K. R., Gennings, C. Jr., Chinchilli, M. V., and Dawson, S. K. (1991). A simple approach for finding estimable functions in linear models, *The American Statistician*, **45**, 51–53.
- Graybill, F. A. (1976). *Theory and Application of the Linear Model*, Wadsworth, California.
- Hartley, H. O. (1967). Expectations, variances and covariances of ANOVA means squares by “synthesis”, *Biometrics*, **23**, 105–114.
- Henderson, C. R. (1953). Estimation of variance and covariance components, *Biometrics*, **9**, 226–252.
- Hicks, C. R. (1973). *Fundamental Concepts in the Design of Experiments*, Holt, Rinehart and Winston, New York.
- Milliken, G. A. and Johnson, D. E. (1984). *Analysis of Messy Data*, Van Nostrand Reinhold, New York.
- Satterthwaite, F. E. (1946). An approximate distribution of estimates of variance components, *Biometrics Bulletin*, **2**, 110–114.
- Searle, S. R., Casella, G., and McCulloch, C. E. (1971). *Linear Models*, John Wiley and Sons, New York.
- Searle, S. R., Casella, G., and McCulloch, C. E. (1992). *Variance Components*, John Wiley and Sons, New York.

혼합모형의 추정가능함수

최재성^{a,1}

^a계명대학교 통계학과

(2015년 10월 27일 접수, 2015년 12월 22일 수정, 2016년 1월 19일 채택)

요약

본 논문은 고정요인과 확률요인의 혼합모형에서 추정가능함수를 논의하고 있다. 고정효과모형에서 정의된 추정가능함수가 혼합효과모형에서 어떻게 정의되어야 하는가를 규정하고 추정가능함수의 분산추정치를 구하는 방법을 제시하고 있다. 또한 혼합모형에서 분산성분의 추정을 위한 제곱합의 계산에 상수적합법을 이용하고 추론을 위한 자유도의 계산에 Satterthwaite의 근사화를 다루고 있으며 분산성분을 구하기 위한 모형의 적합방식으로 단계별 방법을 적용하고 있다. 모형의 단계별 적합에서 주어지는 모형행렬의 사영을 이용한 제1종 제곱합의 계산방식이 제공되며 사영을 이용한 변동요인별 제1종 제곱합의 기댓값 계산에 Hartley의 합성법이 논의된다.

주요용어: 혼합모형, 추정가능함수, 상수적합법, 제1종 제곱합, 합성법

¹(42601) 대구광역시 달서구 달구벌대로 1095, 계명대학교 통계학과. E-mail: jschoi@kmu.ac.kr