

표본조사의 가중치 작업에 대한 고찰

- Composite calibration estimator의 소개 -

김 재 광*

1. 서론

표본조사에서 사용되는 대부분의 추정량은 표본 관측치의 가중 합으로 종종 표현되는데 이렇게 관측치의 일차 가중합으로 추정량을 구현하는 방법은 하나의 가중치가 여러 개의 항목에 공통적으로 쓰이게 됨으로써 다목적 조사(multi-purpose survey)의 추론에 편리하다. 또한, 예를 들어 총수입과 같은 항목은 여러 가지 세부 수입 항목들의 합으로 표현되는데 이렇게 세부 수입 항목치에 대한 통계치 합이 총수입 항목의 통계치와 같아 집으로써 일관성있는 통계치가 구현되기 위해서는 그 통계 추정을 일차 가중합으로 구현해야 할 것이다. 이러한 성질을 통계량의 내적 일관성(internal consistency)이라고도 하는데 이는 누구나 같은 결과를 얻을수 있다는 점에서 특히 자료가 일반에게 공개되는 주요 국가 통계에서 반드시 갖추어야할 중요 성질이라 할 수 있을 것이다.

가중치 작업은 표본 추출을 통해서 얻어지는 표본에 대하여 적절한 가중치를 부여해 줌으로써 추정량의 효율성과 신뢰성을 제고하고자 하는 방법으로 표본 조사에서의 필수적인 요소이다. 가중치 작업은 표본 추출 확률을 이용한 설계 가중치 (design weight)외에도 외부 보조 변수와의 일치성을 위해 조정해 주는 calibration, 단위 무응답 자료의 처리를 위한 무응답 가중치 조정, 그리고 지나치게 가중치 이상값(outlier)를 판별하여 이를 처리하는 가중치 이상값 처리, 이러한 요소들로 구성되어 있다. 이를 정리하면 다음과 같다.

설계 가중치 작성 -> 무응답 가중치 조정 -> calibration -> 가중치 이상값 처리

설계 가중치는 일차 표본 포함 확률(first-order inclusion probability)의 역수를 사용하여 계

*연세대학교 응용통계학과

의 형태로 표현될수 있을 것이다. 이 단순 회귀 추정량은 x 에 대하여 같은 추정량 \bar{x}_1 을 구현해 낸다는 일종의 calibration 성질을 지니는 장점이 있지만 전통적인 회귀 추정량과는 달리 (6)의 추정량의 분산이 y 의 단순 추정량 \bar{y}_2 의 분산보다 더 작게 된다는 보장이 없게 된다.

이러한 문제점을 해결하기 위하여 Zieschang (1990)과 Renssen and Nieuwenbroek (1997)은 \bar{x}_1 을 control로 사용하는 대신

$$\bar{x}_\alpha = \alpha \bar{x}_1 + (1 - \alpha) \bar{x}_2 \quad \dots \quad (7)$$

을 control로 사용하는 regression 추정량을 제안하였다. 즉, (7)을 사용한 회귀 추정량은

$$\widehat{\theta}_{reg} = \bar{y}_2 + (\bar{x}_\alpha - \bar{x}_2)b \quad \dots \quad (8)$$

의 형태로 표현되며 $\alpha = 0$ 인 경우에는 (8)의 회귀 추정량이 단순 추정량 \bar{y}_2 와 동일해지고 $\alpha = 1$ 인 경우에는 (6)의 단순 회귀 추정량이 된다. 최적 계수 α 는 (8)의 분산을 최소화하도록 결정하면 된다.

(8)의 추정량은 또한 다음과 같이 표현된다.

$$\widehat{\theta}_{reg} = \alpha \widehat{\theta}_{sr} + (1 - \alpha) \bar{y}_2 \quad \dots \quad (9)$$

즉, 식 (9)는 식(6)의 단순 회귀 추정량과 단순 추정량 \bar{y}_2 의 가중평균인 일종의 복합 추정량 (composite estimator)의 형태가 된다. 이 복합 추정량은 x 의 추정에 대해 calibration을 유지하므로 최적 계수 α 를 사용한 복합 추정량은 calibration을 유지하면서 분산을 최소화하는 추정량이 될 것이다. 여기서 α 는 $\widehat{\theta}_{sr}$ 을 \bar{y}_2 방향으로 값을 보정해주는 Shrinkage 계수로도 불리울 수 있다. 이 복합추정량의 분산을 최소화하도록 그 계수를 구하면

$$\alpha^* = \frac{Var(\widehat{\theta}_{sr}) - Cov(\bar{y}_2, \widehat{\theta}_{sr})}{Var(\bar{y}_2) + Var(\widehat{\theta}_{sr}) - 2Cov(\bar{y}_2, \widehat{\theta}_{sr})}$$

으로 표현된다.

4. 복합 추정량의 확장

이 절에서는 3절에서 다른 복합 추정량처럼 calibration을 유지하면서 분산을 최소화하는 추정량을 고려해 보기로 한다. 본 절에서는 보다 일반적인 형태의 calibration 조건을 이용한 복합 추정량을 연구하고자 한다.

3절에서처럼 하나의 모두에 대해 두가지 다른 불편 추정량 \bar{x}_1 과 \bar{x}_2 이 존재한다고 하면

참고 문헌

- Deville and Sarndal (1992) "Calibration estimators in survey sampling", *Journal of the American Statistical Association*, 87, p376-382.
- Renssen and Nieuwenbroek (1997) "Aligning estimates for common variables in two or more sample surveys", *Journal of the American Statistical Association*, 92, p368-374.
- Sarndal and Lundstrom (2004) "Estimation in surveys with nonresponse", *Wiley*.
- Zieschang (1990) "Sample weighting methods and estimation of totals in the consumer expenditure survey", 85, p986-1001.