

## بررسی مقایسه ای مدل پاسخ تصادفی بر روی سوال حساس بر اساس رگرسیون لجستیک

علی شمس پور ، دانشجوی کارشناسی ارشد گروه آمار، دانشگاه شهید چمران اهواز، ali\_shamspour@yahoo.com  
دکتر سید محمد رضا علوی ، عضو هیأت علمی، گروه آمار، دانشگاه شهید چمران اهواز alavi\_m@mscu.ac.ir  
دکتر محمدرضا آخوند ، عضو هیأت علمی ، گروه آمار، دانشگاه شهید چمران اهواز mr.akhoond@scu.ac.ir

**چکیده:** مدل پاسخ تصادفی ( $RM$ )<sup>۱</sup> معرفی شده توسط وارنر (۱۹۶۵) ، برای جلوگیری از پاسخ نادرست به سوالات حساس و محافظت از حریم خصوصی مصاحبه شونده طراحی شد. در این مقاله یک رویکرد از مدل نمونه مورد بررسی قرار می گیرد و برآوردی از کل افراد با مشخصه حساس را ارائه می دهد. همچنین مدل های پاسخ تصادفی، مدل رگرسیون لجستیک برای داده های پاسخ تصادفی شده و برآورد  $\pi$  مورد مطالعه قرار می گیرد و با استفاده از شبیه سازی کارایی آن مقایسه می شود.

**کلمات کلیدی:** پاسخ تصادفی ، برآورد رگرسیون لجستیک ، برآورد  $\pi$ .

### مقدمه

این امکان را می دهد، تا علاوه بر امنیت پاسخ دهنده، برآوردی با کارایی مناسب تر در پاسخ به سوال حساس مطرح شود. این روش ابتدا توسط وارنر (۱۹۶۵) برای پاسخ به سوال حساس طراحی شد، به منظور بهبود این روش هارویتس و همکاران (۱۹۶۷) ؛ گرینبرگ و همکاران (۱۹۶۹) طراحی سوال نامربوط را نیز ارائه دادند، در همین حوزه به مدل های پاسخ تصادفی دیور (۱۹۷۶) و مگنت و سینگ (۱۹۹۰) می توان اشاره کرد. همچنین در همین راستا لهتونن-ویجین (۱۹۹۸) ، با استفاده از مدل رگرسیون لجستیک به منظور بهبود کارایی روش پاسخ تصادفی به تحقیق

در مبحث نمونه گیری پاسخ درست به سوال حساس دارای اهمیت زیادی می باشد. زیرا بسیاری از افراد در هنگام مواجهه با سوال حساس، پاسخ واقعی نمی دهند و یا بطور کلی از پاسخ دادن به سوال امتناع می کنند. از جمله مسائل حساسی مانند: استفاده از مواد مخدر، ترجیحات جنسی، مالیات، مشروبات الکلی و... این موضوع باعث می شود که، نتایج بدست آمده با استفاده از روش مستقیم در پاسخ به سوال حساس، دارای اریبی باشند. روش پاسخ تصادفی

<sup>۱</sup> Randomized response technique



دو سوال با محوریت حساس و غیر حساس مطرح می شود. در این مقاله مدل رگرسیون لجستیک ارائه خواهد شد و فرض بر این خواهد بود، که متغیر کمکی حساس نمی باشد. سپس برآورد واریانس برای مدل های مورد نظر ارائه خواهد شد و با استفاده از شبیه سازی کارایی مدل ها مورد بررسی قرار می گیرد.

## چارچوب

جامعه محدود،  $U = \{1, 2, \dots, N\}$  در نظر گرفته می شود.  $y$  متغیر دوحالتی می باشد، که اشاره به این دارد، فرد پاسخ دهنده به کدام گروه تعلق دارد (گروه حساس  $A$  یا گروه غیر حساس  $A^c$ )  $y_k$  مقداری از  $y$  برای  $k^{th}$  عنصر از جامعه می باشد.  $y_k = 1$ ، فرد  $k$ ام پاسخ دهنده متعلق به ویژگی حساس  $A$  می باشد و  $y_k = 0$ ، فرد  $k$ ام پاسخ دهنده متعلق به ویژگی غیر حساس  $A^c$  می باشد. هدف بدست آوردن برآورد  $t_A = \sum_U y_k$  می باشد.  $t_A$  تعداد کل افراد با مشخصه حساس می باشد. روند نمونه گیری شامل دو مرحله می باشد: مرحله اول شامل انتخاب نمونه، نمونه ای به اندازه  $n$ ، از جامعه مورد نظریا توجه به طرح نمونه گیری  $p(s)$  با احتمال شمول مثبت و مرحله دوم مصاحبه از افراد نمونه مطابق با مدل تصادفی صورت می گیرد و برای هر  $k \in s$ ، مدل پاسخ تصادفی را با متغیر  $z_k$  معرفی می کنیم. چنان که:  $\hat{z}_k = az_k + b_k$  برآوردی نا اریبی برای  $y_k$  می باشد.  $a, b_k$  ثابت های شناخته شده می باشند.

## برآورد رگرسیون لجستیک بر روی مدل های پاسخ تصادفی

$t_{AG,LRGEG}$  برآورد رگرسیون لجستیک تعمیم یافته  $t_A$  می باشد. فرض می شود:  $\underline{y} = \{y_1, y_2, \dots, y_N\}$  برگرفته شده از جامعه  $\underline{Y} = \{Y_1, Y_2, \dots, Y_N\}$

<sup>†</sup> Generalized Logistic Regression Model Estimator

پرداختند. مدل  $W$  (وارنر ۱۹۶۵) این مدل با استفاده از ابزار تصادفی دو سوال با یک پاسخ بله و خیر را مطرح می کند، سوال اول: آیا شما متعلق به گروهی با ویژگی  $A$  می باشید؟ سوال دوم: آیا شما متعلق به گروهی با ویژگی  $A^c$  می باشید؟  $A$  مشخصه حساس مورد توجه می باشد. به عنوان مثال: سوال اول اینگونه مطرح می شود: آیا شما مالیات سال گذشته خود را پرداخت کرده اید؟ سوال دوم: آیا مامور مالیات برای ثبت مشخصات شما در سال گذشته مراجعه کرده است؟ مدل  $U$  (گرینبرگ و همکاران ۱۹۶۹) مدل پاسخ تصادفی در طراحی سوال حساس نامربوط می باشد و همانند مدل  $W$ ، با استفاده از ابزار تصادفی دو سوال با یک پاسخ بله و خیر مطرح می کند، تفاوت این روش با وارنر نامربوط بودن سوال غیر حساس نسبت به سوال حساس می باشد. برای مثال: آیا شما در سال گذشته مالیات خود را پرداخت کرده اید؟ سوال دوم: آیا شما علاقمند به تماشای سریال های تلویزیونی می باشید؟ مدل  $D$  (مدل دیور و همکاران ۱۹۷۷) همانند مدل  $U$  می باشد، با این تفاوت که پاسخ مربوط به سوال غیر حساس نامربوط، پاسخی کاملا مشخص می باشد. مدل  $H$  (هارویتس و همکاران ۱۹۷۶) این مدل برای طراحی سوال نامربوط می باشد، در این مدل سه سوال مطرح می شود: سوال اول، یک سوال حساس سوال دوم، سوالی غیر حساس که پاسخ آن مثبت می باشد و سوال سوم، سوالی غیر حساس که پاسخ آن منفی می باشد. که به ترتیب احتمال انتخاب سوالات  $P_1, P_2, P_3$  می باشد. بطوری که  $P_1 + P_2 + P_3 = 1$ . مدل  $M$  (مگنت و سینگ ۱۹۹۰) این مدل از دو ابزار تصادفی استفاده می کند، ابتدا دو سوال مطرح می شود: سوال اول، سوالی حساس و سوال دوم سوال غیر حساس نامربوط می باشد، در صورتی که پاسخ مصاحبه کننده، مربوط به سوال نامربوط باشد، مصاحبه کننده وارد ابزار تصادفی دوم می شود، که در اینجا نیز

نشانگر شمول: فرض کنید  $k \in s$ ، نشان دهنده این پیشامد باشد، که نمونه  $s$  واحد  $i$  را در بر داشته باشد. متغیر تصادفی

$$I_k(s) = \begin{cases} 1 & 1 \leq k \leq N, K \in s \\ 0 & \text{otherwise} \end{cases}$$

نشانگر شمول معرفی می کنیم. مدل (W):

$$Z_k = \begin{cases} y_k & p \\ 1 - y_k & 1 - p \end{cases}$$

به ازای هر  $k \in s$

$$\hat{z}_k = \frac{z_k - (1 - p)}{2p - 1}$$

مدل (D):

$$Z_k = \begin{cases} y_k & p \\ 1 & 1 - p \end{cases}$$

$$E_R = (Z_K) = y_k p + 1 - p$$

$$\hat{z}_k = \frac{z_k - (1 - p)}{p}$$

مدل (U):

$$Z_k = \begin{cases} y_k & p \\ w_k & 1 - p \end{cases}$$

$$E_R = (Z_K) = y_k p + w_k (1 - p)$$

$$\hat{z}_k = \frac{z_k - w_k (1 - p)}{p}$$

<sup>۳</sup> Estimator of The Variance Estimator

می باشد. این بردار مشاهدات متغیر تصادفی مستقل هستند و همچنین دارای توزیعی می باشند که در زیر معرفی می کنیم:

$$P(Y_k = 1 | x_k, \beta) = \frac{\exp(x_k, \beta)}{1 + \exp(x_k, \beta)}$$

$$k = \{1, 2, \dots, N\}$$

$$t_A = \sum_U Y_K = \sum_U \mu_k + \sum_U (Y_K - \mu_k)$$

معادله زیر برآورد (t<sub>A</sub>) را معرفی می کند:

$$(\hat{t}_A)_{G,L RGE G} = \sum_U \hat{\mu}_k + \sum_s \frac{\hat{z}_k - \hat{\mu}_k}{\hat{\mu}_k}$$

$$\hat{\mu}_k = \mu\{x_k^t, \hat{\beta}\} = \frac{1}{1 + \exp(-x_k^t, \hat{\beta})}$$

$\hat{\beta}$  برآورد درست‌نمایی ماکزیمم  $\beta$  می باشد. که با الگوریتم نیوتن رافسون محاسبه می شود. همچنین در زیر معرفی می کنیم: برآورد  $\pi$ ، برآوردی از واریانس تعداد کل افرادی که به سوال حساس پاسخ می دهند. <sup>۳</sup> (سارندال و همکاران ۱۹۹۲؛ ولتن و ویجانن ۱۹۹۸)

$$\hat{V}(\hat{t}_{AG,L RGE G}) = a^2 \sum_s \sum_{kl} \frac{\Delta_{kl}}{\pi_{kl}} \left( \frac{z_k - \hat{\lambda}_k}{\pi_k} \right) \left( \frac{z_l - \hat{\lambda}_l}{\pi_l} \right)$$

$$= \sum_s \sum_{kl} \frac{\Delta_{kl}}{\pi_{kl}} \left( \frac{\hat{z}_k - \hat{\lambda}_k}{\pi_k} \right) \left( \frac{\hat{z}_l - \hat{\lambda}_l}{\pi_l} \right)$$

$$cov(I_K, I_l) = \Delta_{kl} = \pi_{kl} - \pi_k \pi_l$$

$$\hat{\lambda}_k = \frac{\hat{\mu}_k - b_k}{a}$$

$$\hat{\lambda}_l = \frac{\hat{\mu}_l - b_l}{a}$$

$$k, l = 1, 2, \dots, N$$

## پیوستار

طرح نمونه گیری با احتمال شمول مثبت  $\pi_{kl} \leq \pi_k$  در زیر معرفی می کنیم:

$$\pi_k = \sum_{s \in s} p(s), \pi_{kl} = \sum_{k, l \in s} p(s)$$



## نتایج

توجهی دربرآورد واریانس دارد.

## مراجع

- [1] Soberanis-Cruz, Víctor Hugo, and Víctor Miranda-Soberanis. "The Generalized Logistic Regression Estimator in a Finite Population Sampling without Replacement Setting with Randomized Response." *Revista Colombiana de Estadística* 34.3 (2011): 451-46 *Envirometrics*.
- [2] LEHTONEN, R., and VEIJANEN, A., (1998): *Logistic Generalized Regression Estimators. Survey*. New York: Springer Verlag.
- [3] SÄRNDAL, C.E., SWENSSON, B., and WRETMAN, J. (1992): *Model Assisted Survey Sampling*. New York: Springer Verlag.
- [4] HORVITZ, D.C., GREENBERG, B. G., and ABERNATHY, J. R. (1976): RR: a data gathering device for sensitive questions. *Internat. Statist. Rev.* 44, 181-196
- [5] Devore, Jay L. "A note on the randomized response technique." *Communications in Statistics-Theory and Methods* 6.15 (1977): 1525-1529. *Methodology*, 24, 51-55.
- [6] Warner, Stanley L. "Randomized response: A survey technique for eliminating evasive answer bias." *Journal of the American Statistical Association* 60.309 (1965): 63-69.
- [7] Mangat, N. S., and Ravindra Singh. "An alternative randomized response procedure." *Biometrika* 77.2 (1990): 439-442.

در این مقاله سه روش  $W$ ،  $U$  و  $D$  از نظر کارایی مورد مقایسه قرار می گیرد و نتایج شبیه سازی به شرح زیر خواهد بود.

	$\hat{t}_A$	$\sqrt{\text{var}(\hat{t}_A)}$	$\sqrt{\widehat{\text{var}}(\hat{t}_A)}$
W	۴۷۳,۲۳۱	۷۶,۰۲۵	۶۵,۴۶
D	۴۷۵,۸۸	۷۶,۰۲	۳۱,۴۴

با توجه به جدول بالا که بررسی مقایسه ای دو مدل وارنر و دیور می باشد، و همان طور که ملاحظه می کنید: مدل  $D$  نسبت به مدل  $W$  برآوردی با کارایی مناسب تری برای تعداد کل افراد جامعه با صفت حساس دارد. و دارای ارزیابی کمتری نسبت به مدل وارنر می باشد. در ادامه مقایسه مدل وارنر و مدل گرینبرگ و همکاران مطرح شده است و همانطور که در جدول زیر ملاحظه می کنید: مدل  $U$  نسبت به مدل  $W$  برآوردی با کارایی مناسبتری برای تعداد کل افراد جامعه با صفت حساس دارد و دارای ارزیابی کمتری نسبت به مدل وارنر می باشد.

	$\hat{t}_A$	$\sqrt{\text{var}(\hat{t}_A)}$	$\sqrt{\widehat{\text{var}}(\hat{t}_A)}$
W	۴۷۳,۲۳۱	۷۶,۰۲۵	۶۵,۴۶
U	۴۷۴,۸۵	۷۴,۱۴۴	۳۰,۸۳

همچنین از مقایسه دو جدول بالا به این نتیجه نیز خواهیم رسید که، مدل  $U$  نسبت به مدل  $D$  دارای ارزیابی کمتری خواهد بود. در این نتیجه گیری به این موضوع پی خواهیم برد، که پاسخ به سوال حساس با پرسش مستقیم با مشکل جدی مواجه هست، این مقاله برای بررسی این موضوع، به برآورد رگرسیون لجستیک با توجه به مدل پاسخ تصادفی اشاره کرده است و همچنین، روش مناسبی برای برآورد تعداد کل افراد با مشخصه حساس را نشان می دهد، که نتایج حاکی از این است که مدل  $U$  دارای کارایی بهتری نسبت به دو روش ذکر شده می باشد و این روش کاهش قابل