

## افزایش دقت فواصل اطمینان بوت‌استرپی برای خط رگرسیون

محسن طیبیان، سیدمحمدابراهیم حسینی نسب  
گروه آمار، دانشگاه تربیت مدرس  
گروه آمار، دانشگاه شهید بهشتی

در این مقاله روش‌های بوت‌استرپ را برای تشکیل فاصله اطمینان به روش  $t$ -صدکی برای خط رگرسیون استفاده می‌کنیم. هدف ما از این مقاله، بهبود مرتبه دقت فواصل اطمینان می‌باشد که برای این منظور، احتمال پوشش و خطای پوشش این فاصله را محاسبه می‌کنیم که خطای پوشش برای فاصله بوت‌استرپی منفرد به روش  $t$ -صدکی به صورت  $O(n^{-1})$  و برای فاصله بوت‌استرپی دوگانه به روش  $t$ -صدکی برابر  $O(n^{-2})$  می‌باشد.

واژه‌های کلیدی: احتمال پوشش، خطای پوشش، فاصله اطمینان بوت‌استرپی.

### ۱ مقدمه

اولین استفاده از بوت‌استرپ برای برآورد توزیع برآوردهای ضرایب رگرسیون خطی توسط افرون (۱۹۷۹) مطرح شد، و توسط فریدمن (۱۹۸۱) ادامه یافت. این روش بر تقریب توزیع خطاهای مشاهده نشده با توزیع تجربی مانده‌های مرکزی شده استوار است. نویدی (۱۹۸۹) نشان داد زمانی که عضوهای ماتریس طرح بدون محدودیت افزایش یابند، بوت‌استرپ به طور مجانبی بهتر از تقریب نرمال است. وی همچنین مطرح کرد که در موارد دیگر، همواره بوت‌استرپ به خوبی تقریب نرمال می‌باشد. هال (۱۹۹۰) و چرنیک (۲۰۰۸) نشان دادند زمانی که توزیع مانده‌ها در مدل نرمال نباشد و به خصوص اینکه توزیع آنها دارای دم‌های کلفت باشد می‌توان از روش‌های بوت‌استرپ در رگرسیون استفاده کرد. مارتین و رابرتز (۲۰۰۶) و چرنیک (۲۰۰۸) نیز نشان دادند اگر تعدادی مشاهده‌ی پرت وجود در مدل داشته باشند آنگاه می‌توان از روش‌های بوت‌استرپ استفاده کرد. فیشر و هال (۱۹۹۱) نشان دادند که روش بوت‌استرپ برای نمونه‌های با حجم کوچک، عملکرد مناسبی دارد. همچنین هال (۱۹۹۲) فواصل اطمینان بوت‌استرپی را به روش‌های صدکی و  $t$ -صدکی برای پارامتر شیب در مدل رگرسیونی خطی ساده بدست آورد و احتمال و خطای پوشش این فواصل را

بررسی کرد.

در این مقاله ابتدا، فواصل اطمینان بوت استرپی را به روش  $t$ -صدکی برای خط رگرسیونی بدست می آوریم و سپس احتمال پوشش و خطای پوشش این فواصل را نیز محاسبه می کنیم. در انتها دقت فواصل اطمینان بوت استرپی منفرد و دوگانه را به طور نظری بدست آورده و روش های عددی برای تشکیل فواصل اطمینان بوت استرپی منفرد و دوگانه مطرح می کنیم.

## ۲ مدل خطی ساده

مدل خطی ساده زیر را در نظر بگیرید:

$$Y_i = \beta_0 + X_i \beta_1 + \epsilon_i, \quad i = 1, \dots, n,$$

که در آن  $X_i$  ها نقاط طرح اند و ثابت فرض می شوند و جملات خطا،  $\epsilon_i$  ها، مستقل و هم توزیع با میانگین صفر و واریانس  $\sigma^2$  هستند. برآوردهای  $\hat{\beta}_0$ ،  $\hat{\beta}_1$  و  $\hat{\sigma}^2$  به ترتیب برای پارامترهای  $\beta_0$ ،  $\beta_1$  و  $\sigma^2$  به صورت زیر تعریف می شوند:

$$\hat{\beta}_0 = \bar{Y} - \bar{X} \hat{\beta}_1, \quad \hat{\beta}_1 = n^{-1} \sigma_x^{-2} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}), \quad \hat{\sigma}^2 = n^{-1} \sum_{i=1}^n \hat{\epsilon}_i^2, \quad (1)$$

که در آن داریم:

$$\sigma_x^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

فرض کنید چولگی و کشیدگی جامعه به ترتیب به صورت زیر باشند:

$$\gamma = E(\epsilon/\sigma)^3, \quad \kappa = E(\epsilon/\sigma)^4 - 3. \quad (2)$$

همچنین فرض کنید که برآورد تجربی آنها را با روابط زیر نشان دهیم:

$$\hat{\gamma} = n^{-1} \sum_{i=1}^n (\hat{\epsilon}_i/\hat{\sigma})^3, \quad \hat{\kappa} = n^{-1} \sum_{i=1}^n (\hat{\epsilon}_i/\hat{\sigma})^4 - 3.$$

با استفاده از روش باز نمونه گیری از مانده ها در رگرسیون داریم:

$$Y_i^* = \hat{\beta}_0 + X_i \hat{\beta}_1 + \epsilon_i^*, \quad 1 \leq i \leq n,$$

که در آن  $\epsilon_i^*$ ها به طور تصادفی و با جایگذاری از مانده‌های  $\epsilon_i$  استخراج شده‌اند. به علاوه برای محاسبه‌ی  $\hat{\beta}_0^*$ ،  $\hat{\beta}_1^*$  و  $\hat{\sigma}^*$  کافی است که به ترتیب در فرمول‌های  $\hat{\beta}_0$ ،  $\hat{\beta}_1$  و  $\hat{\sigma}$ ، که در (۱) معرفی شده‌اند به جای استفاده از  $Y_i$ ، از مقدار  $Y_i^*$  استفاده شود. برآورد مقدار میانگین

$$d_0 = E(Y|X = x_0) = \beta_0 + x_0\beta_1,$$

عبارت است از:

$$\hat{d}_0 = \hat{\beta}_0 + x_0\hat{\beta}_1. \quad (3)$$

اگر قرار دهیم  $x_0 = 0$ ، آنگاه  $\hat{\beta}_0$  برآورد عرض از مبدا  $\beta_0$  بدست می‌آید. بنابراین مساله‌ی برآورد  $d_0$ ، یک مساله کلی است. همچنین داریم:

$$Var(\hat{d}_0) = n^{-1}\sigma_d^2,$$

که در آن

$$\sigma_d^2 = 1 + \sigma_x^{-2}(x_0 - \bar{x})^2. \quad (4)$$

در نتیجه آماره زیر (به طور مجانبی) یک کمیت محوری است:

$$T = n^{1/2}(\hat{d}_0 - d_0)/(\hat{\sigma}\sigma_d). \quad (5)$$

با روش بازنمونه‌گیری از مانده‌ها در مدل‌های رگرسیونی، بازنمونه‌ی  $\{(X_1, Y_1^*), (X_2, Y_2^*), \dots, (X_n, Y_n^*)\}$  را تولید می‌کنیم. فرض کنید که  $\hat{\beta}_0^*$ ،  $\hat{\beta}_1^*$  و  $\hat{\sigma}^{*2}$  به ترتیب دارای فرمول یکسانی با  $\hat{\beta}_0$ ،  $\hat{\beta}_1$  و  $\hat{\sigma}^2$  که در (۱) معرفی شده‌اند باشند به جز اینکه برای محاسبه‌ی آنها به جای استفاده از  $Y_i$  از مقدار  $Y_i^*$  استفاده شود. نسخه‌های بوت‌استرپی  $\hat{d}_0^*$  و  $T^*$  به ترتیب عبارتند از:

$$\begin{aligned} \hat{d}_0^* &= \hat{\beta}_0^* + x_0\hat{\beta}_1^*, \\ T^* &= n^{1/2}(\hat{d}_0^* - \hat{d}_0)/(\hat{\sigma}^*\sigma_d). \end{aligned} \quad (6)$$

### ۳ دقت فواصل اطمینان بوت‌استری برای خط رگرسیون

فاصله اطمینان دو طرفه با دم‌های برابر به روش  $t$ -صدکی برای خط رگرسیون  $(d_0)$  با سطح پوشش اسمی  $1 - \alpha$  عبارت است از:

$$J_T = \left( \hat{d}_0 - n^{-1/2} \hat{\sigma}_{\sigma_d} v_{1-\alpha/2}, \hat{d}_0 - n^{-1/2} \hat{\sigma}_{\sigma_d} v_{\alpha/2} \right),$$

که در آن مقادیر  $\hat{\sigma}$ ،  $\hat{d}_0$  و  $\sigma_d$  به ترتیب در روابط (۱)، (۳) و (۴) معرفی شده‌اند. چندک‌های  $v_{1-\alpha/2}$  و  $v_{\alpha/2}$  به ترتیب جواب‌های معادلات  $P(T \leq v_{1-\alpha/2}) = 1 - \alpha/2$  و  $P(T \leq v_{\alpha/2}) = \alpha/2$  هستند و آماره  $T$  در رابطه (۵) معرفی شده است. نسخه بوت‌استری این فاصله به صورت زیر می‌باشد:

$$\hat{J}_T = \left( \hat{d}_0 - n^{-1/2} \hat{\sigma}_{\sigma_d} \hat{v}_{1-\alpha/2}, \hat{d}_0 - n^{-1/2} \hat{\sigma}_{\sigma_d} \hat{v}_{\alpha/2} \right), \quad (7)$$

که در آن چندک‌های  $\hat{v}_{1-\alpha/2}$  و  $\hat{v}_{\alpha/2}$  به ترتیب جواب‌های معادلات  $P(T^* \leq \hat{v}_{1-\alpha/2} | \chi) = 1 - \alpha/2$  و  $P(T^* \leq \hat{v}_{\alpha/2} | \chi) = \alpha/2$  هستند و آماره  $T^*$  در رابطه (۶) معرفی شده است. لازم بذکر است که فواصل  $J_T$  و  $\hat{J}_T$ ، به هر دم احتمال مساوی تخصیص می‌دهند، یعنی برای مثال برای فاصله  $J_T$  داریم:

$$P(d_0 \leq \hat{d}_0 - n^{-1/2} \hat{\sigma}_{\sigma_d} v_{1-\alpha/2}) = P(d_0 \geq \hat{d}_0 - n^{-1/2} \hat{\sigma}_{\sigma_d} v_{\alpha/2}) = \alpha/2.$$

**قضیه ۱.** فرض کنید که فاصله اطمینان بوت‌استری با سطح پوشش اسمی  $1 - \alpha$  که از روش  $t$ -صدکی بدست آمده است به صورت فاصله  $\hat{J}_T$  باشد که در رابطه‌ی (۷) معرفی شده است، آنگاه احتمال پوشش آن برابر است با:

$$P(d_0 \in \hat{J}_T) = 1 - \alpha - \frac{n^{-1}}{\gamma} \sigma_d^{-1} (\kappa - \frac{3}{\gamma} \gamma^2) z_{1-\alpha/2} \\ \times \left\{ (\gamma d - 3 \sigma_z^{-1}) d_{1-\alpha/2}^2 + \gamma d \right\} \phi(z_{1-\alpha/2}) + O(n^{-2}),$$

که در آن مقادیر  $\gamma$ ،  $\kappa$  و  $\sigma_d$  در روابط (۲) و (۴) معرفی شده‌اند و  $\phi(\cdot)$  تابع چگالی نرمال استاندارد می‌باشد.

اثبات: برای اثبات به حسینی‌نسب و طیبیان (۲۰۰۹) رجوع شود. بنابراین با توجه به قضیه ۱ می‌توان گفت که فاصله اطمینان بوت‌استری  $\hat{J}_T$  دارای خطای پوشش از مرتبه  $O(n^{-1})$  است.



را بدست می آوریم:  $\hat{v}_{b,\alpha}^* = \hat{v}_{b,\alpha/2}^*, \hat{v}_{b,1-\alpha/2}^*$

$$\frac{1}{C} \sum_{c=1}^C I\left(\frac{n^{1/2}(\hat{d}_{o,bc}^{**} - \hat{d}_{o,b}^*)}{\sigma_d \hat{\sigma}^{**}} \leq \hat{v}_{b,\alpha}^*\right) = \alpha,$$

که در آن مقادیر  $\hat{d}_{o,bc}^{**}$  و  $\hat{\sigma}^{**}$  از بازنمونه  $\chi_{bc}^{**} = \{(X_1, Y_{1bc}^{**}), \dots, (X_n, Y_{nbc}^{**})\}$  محاسبه شده‌اند. حال از حل معادلات زیر،  $\tilde{v} = \tilde{v}_{\alpha/2}, \tilde{v}_{1-\alpha/2}$  بدست می آید:

$$\frac{1}{B} \sum_{b=1}^B I\left(\frac{n^{1/2}(\hat{d}_{o,b}^* - \hat{d}_o)}{\sigma_d \hat{\sigma}^*} \leq \hat{v}_{b,\alpha/2}^* + \tilde{v}\right) = \alpha/2,$$

$$\frac{1}{B} \sum_{b=1}^B I\left(\frac{n^{1/2}(\hat{d}_{o,b}^* - \hat{d}_o)}{\sigma_d \hat{\sigma}^*} \leq \hat{v}_{b,1-\alpha/2}^* + \tilde{v}\right) = 1 - \alpha/2.$$

بنابراین فاصله اطمینان بوت استرپ دوگانه با دم‌های برابر برای خط رگرسیونی با سطح پوشش اسمی  $1 - \alpha$  عبارت است از:

$$\hat{J}_{\tau d} = \left(\hat{d}_o - n^{-1/2} \sigma_d \hat{\sigma} (\hat{v}_{1-\alpha/2} + \tilde{v}_{1-\alpha/2}), \hat{d}_o - n^{-1/2} \sigma_d \hat{\sigma} (\hat{v}_{\alpha/2} + \tilde{v}_{\alpha/2})\right).$$

**قضیه ۲.** اگر فاصله اطمینان بوت استرپی دوگانه به روش  $t$ -صدکی برای خط رگرسیون مانند بالا بدست آید، آنگاه:

$$P(d_o \in \hat{J}_{\tau d}) = 1 - \alpha + O(n^{-2}),$$

که در آن خطای پوشش به صورت  $O(n^{-2})$  می‌باشد. یعنی، در صورت استفاده از روش بوت استرپ دوگانه، خطای پوشش فاصله از  $O(n^{-1})$  به  $O(n^{-2})$  کاهش می‌یابد.

اثبات: برای اثبات به حسینی‌نسب و طبیبیان (۲۰۰۹) رجوع شود.

## ۵ نتیجه‌گیری

خطای پوشش برای فاصله اطمینان بوت استرپی منفرد به صورت  $O(n^{-1})$  می‌باشد در حالیکه خطای پوشش فاصله اطمینان بوت استرپی دوگانه برابر  $O(n^{-2})$  است، یعنی در صورت استفاده از روش بوت استرپ دوگانه، خطای پوشش کاهش می‌یابد.

## مراجع

- [1] Chernick, M. R. (2008). *Bootstrap Methods: A Guide for Practitioners and Researchers*, 2nd ed, United BioSource Corporation, Wiley.
- [2] Efron, B. (1979). Bootstrap methods: Another look at the jackknife. *Ann. Statist.* **7**, 1-26.
- [3] Fisher, N. I. and Hall, P. (1991). Bootstrap algorithms for small samples. *Journal of Statistical Planning and Inference*, **27**, 157-169.
- [4] Freedman, D. (1981). Bootstrapping regression models. *Ann. Statist.*, **9**, 1218-1228.
- [5] Hall, P. (1990). Asymptotic Properties of the Bootstrap for Heavy-Tailed Distributions. *Annals of Probability*, **18**, No. 3, 1342-1360.
- [6] Hall, P. (1992). *The Bootstrap and Edgeworth expansion*, Springer-Verlag, New York.
- [7] Hosseinasab, S.M.E. and Tabibian, M. (2009). Accuracy of Bootstrap Prediction intervals based on simple regression model. Manuscript.
- [8] Martin, M. A. and Roberts, S. (2006). An evaluation of bootstrap methods for outlier detection in least squares regression. *Journal of Applied Statistics*, **33**, No. 7, 703-720.
- [9] Navidi, W. (1989). Edgeworth expansion for Bootstrapping regression models. *Ann. Statist.* **17**, 1472-1478.