

روش های باز نمونه گیری در تحلیل داده های پیچیده

رضا مکرم^۱، وحید رنجبر^۲

چکیده:

در تحلیل داده های پیچیده^۳، معمولاً الگوهای نمونه گیری ساختاری غیر $i.i.d$ را برای داده ها در نظر می گیرند. لذا تکنیک های موجودی که برای برآورد نقطه ای و یا فواصل اطمینان وجود دارند معمولاً برای این نوع از داده ها قابل اجرا نمی باشند. برای حل این مشکل می توان از روش های نمونه گیری مانند نمونه گیری تکراری یا باز نمونه گیری^۴ استفاده کرد که در آنها برآورد پارامتر بدون نیاز به شرط $i.i.d$ بودن داده ها انجام گیرد. از این روش ها می توان با قدرت محاسباتی زیاد همچنین دوری جستن از تئوری های پیچیده برای برآورد پارامتر مورد استفاده قرار گیرند. این روش ها هر کدام دارای محاسن و همچنین کاستی هایی هستند که باعث ایجاد محدودیت در اجرای آنها می شوند. در این مقاله سعی شده است روش های مختلف باز نمونه گیری در الگوی نمونه گیری طبقه ای معرفی و سپس با یک مثال واقعی مورد مقایسه قرار گیرند.

واژه های کلیدی: بوت استرپ با جایگذاری، بوت استرپ بدون جایگذاری، جک نایف، نمونه گیری طبقه ای.

۱ مقدمه

نمود و آن را بوت استرپ نامید. او نشان داد که این روش می تواند از دیگر رقبای خود بهتر باشد اما در این باره بحث های متفاوتی نیز وجود دارد که برای توضیح بیشتر می توان به لو [۷] و وو [۱۳] مراجعه کرد.

هدف ما در این مقاله معرفی بسط روشهای بوت استرپ استاندارد (برای داده های $i.i.d$) به داده های پیچیده است. داده های پیچیده به داده های گفته می شود که برای آنها نمونه گیری به گونه ای است که خاصیت $i.i.d$ وجود ندارد. داده هایی که از روش های نمونه گیری طبقه ای، خوشه ای و نمونه گیری با احتمالات نابرابر

روش های باز نمونه گیری، خصوصاً بوت استرپ و جک نایف، روش هایی هستند که بوسیله آنها می توان به راحتی برآوردگرهایی برای واریانس یافت و همچنین می توان فواصل اطمینان ناپارامتری را برای پارامترهای مورد نظر در جامعه به دست آورد. این روش ها از تئوری بسیار ساده ای نتیجه شده اند و استفاده از آنها غالباً به روش های برنامه نویسی کامپیوتر نیاز دارند. افرن [۴] برای اولین بار روشی جالب از باز نمونه گیری را برای پارامتر مورد نظر $\theta = \theta(F)$ بر اساس یک نمونه تصادفی n تایی $i.i.d$ از جامعه ای با توزیع نامشخص F معرفی

^۱ هیات علمی دانشگاه آزاد اسلامی واحد گرگان، Email: re612m@gmail.com

^۲ استادیار دانشگاه گلستان، گروه آمار، Email: v.ranjbar@gu.c.ir

^۳ Complex data

^۴ Resampling

حاصل می شوند این گونه اند.

روش هایی که در برآورد پارامتر و همچنین فواصل اطمینان برای داده های مختلط وجود دارند در عمل دارای مشکلاتی بوده و اغلب قابل تعمیم به مدل های پیچیده نیستند و همچنین یافتن برآوردگرهای غیر خطی برای آنها مشکل است. سه روشی که عمدتاً برای این داده ها مورد استفاده قرار می گیرند عبارتند از روش خطی سازی^۵ (یا تیلور)، روش جک نایف و روش تکثیرهای متعادل تکرار شده^۶ (*BRR*). کرسکی و راتو [۶] به تفسیر و مقایسه این روش ها پرداخته و برای آنها نتایج مهمی را نیز به دست آوردند.

از مواردی که در بحث برآورد نقطه ای بسیار مورد توجه قرار دارد یافتن برآوردگری برای توابع غیر خطی از میانگین چند متغیر است. از جمله مهمترین این توابع ضریب همبستگی، نسبت و یا رگرسیون چند متغیر است. ممکن است یافتن این برآوردگرها در برخی مواقع چندان مشکل نباشد، اما می دانیم پس از انتخاب برآوردگر مسئله مهم تعیین میزان دقت و یا اصطلاحاً کارایی آن است. کارایی یک برآوردگر را به روش های مختلفی می توان تعیین کرد. از جمله ابزار مورد استفاده در تعیین کارایی، واریانس برآوردگر است. ممکن است یافتن واریانس توابع خطی از میانگین چند متغیر چندان کار سختی نباشد. اما در حالت غیر خطی مطمئناً این گونه نیست. به عنوان مثال فرم بسته ای برای یافتن واریانس ضریب همبستگی دو متغیر وجود ندارد. در این حالت می

توان از روش های باز نمونه گیری استفاده کرد. کرسکی و راتو [۶] نشان دادند که در این سه مدل برای توابع غیر خطی از میانگین چند متغیر مانند رگرسیون یا ضریب همبستگی می توان برآوردگرهایی پیدا کرد که به طور مجانبی سازگار هستند. با این حال این سه روش دارای نقاط ضعفی نیز هستند. به عنوان مثال روش خطی سازی نیازمند محاسبات تئوری زیادی بوده که باعث ایجاد مشکل در اجرای این روش می شود. از طرفی روشهای *BRR* و جک نایف فقط برای نمونه گیری طبقه ای که نمونه گیری اولیه با جایگذاری باشد تعریف شده اند. با این تفصیل ما به روشی نیاز داریم که از محاسبات سخت تئوری دوری کرده و قابل تعمیم به مدل های پیچیده نیز باشد.

در این مقاله بسط روش بوت استرپ به مدل های پیچیده نمونه گیری که در آن نمونه اولیه به صورت بدون جایگذاری استخراج شده است و مشاهدات غیر *i.i.d* هستند را توضیح داده خواهد شد. در قسمت دوم مقاله به نمونه گیری طبقه ای پرداخته می شود و در قسمت سوم روش های بوت استرپ موجود را برای این الگوی معرفی می کنیم که شامل روش بوت استرپ استاندارد (*WR*)، بوت استرپ با جایگذاری^۷ (*BWR*)، روش باز مقیاس گذاری بوت استرپ^۸ (*RS*) و بوت استرپ بدون جایگذاری^۹ (*BWO*) است. در قسمت چهارم مقاله روش *M.M* را که شامل ویژگیهای مثبت روشهای *BWO* و *BWR* و قابل تعمیم به الگوهای

^۵ Linearization

^۶ Balanced Repeated Replication

^۷ With-Replacement Bootstrap

^۸ Rescaling Bootstrap

^۹ Without-Replacement Bootstrap

داده است:

$$Var(\widehat{Y}_{st}) = \sum_{h=1}^L W_h^2 \frac{1-f_h}{n_h} s_h^2. \quad (1)$$

که در آن $f_h = n_h/N_h$ کسر نمونه گیری و s_h^2 برآورد واریانس نمونه ای است، یعنی

$$s_h^2 = \frac{1}{n_h-1} \sum_{i=1}^{n_h} (y_{hi} - \bar{y}_h)^2.$$

هدف در این مقاله یافتن برآوردگری برای رابطه ۱ به وسیله روش های باز نمونه گیری است.

۳ روش های بوت استرپ موجود

یکی از روش هایی که می توان به عنوان یک روش باز نمونه گیری در الگوی نمونه گیری طبقه ای در نظر گرفت اجرای روش بوت استرپ استاندارد برای نمونه در هر طبقه یعنی $\{y_{hi}\}_{i=1}^{n_h}$ است. این روش را با نماد (WR) نمایش می دهیم. در این صورت الگوریتم زیر حاصل می شود.

(۱) نمونه گیری تصادفی با جایگذاری (WR) از نمونه اصلی یعنی $\{y_{hi}\}_{i=1}^{n_h}$ در هر یک از طبقه ها به طور مستقل تا حصول نمونه بوت استرپ یعنی $s^* = \{y_{hi}^* : h = 1, 2, \dots, L; i = 1, 2, \dots, n_h\}$ سپس محاسبه $\hat{\theta}^* = \hat{\theta}(s^*)$.

(۲) B بار تکرار مرحله ۱ و به دست آوردن $\hat{\theta}_B^*, \dots, \hat{\theta}_\tau^*, \hat{\theta}_\lambda^*$

(۳) برآورد $\hat{\theta}$ با مقدار زیر:

$$v_b = E_*(\hat{\theta}^* - E_*\hat{\theta}^*)^2.$$

پیچیده با جایگذاری است را بیان می کنیم. در نهایت در قسمت پنجم با استفاده از یک مثال واقعی از جامعه ای با ساختار طبقه ای، پارامترهای میانگین و واریانس جامعه را توسط روش های بیان شده در این مقاله برآورد کرده و مقایسه می کنیم.

۲ نمونه گیری طبقه ای

در الگوی نمونه گیری طبقه ای، جامعه متشکل از N عضو بوده که به L طبقه مجزا تقسیم شده است. حجم جامعه در هر طبقه به ترتیب برابر با N_1, N_2, \dots, N_L است. در این الگو نمونه گیری به این طریق بوده که از هر طبقه نمونه ای تصادفی ساده بدون جایگذاری (wtr) انتخاب شده که حجم نمونه ها در طبقات به ترتیب برابر است با n_1, n_2, \dots, n_L بطوریکه حجم کل نمونه برابر $n = n_1 + n_2 + \dots + n_L$ خواهد بود.

فرض کنید y_{hi} اندازه ای از یک ویژگی از یک عضو در نمونه باشد که در آن اندیس h بیانگر شماره طبقه و زیرنویس i بیانگر نامین عضو در طبقه است. فرض کنید $\theta = \theta(S)$ پارامتر مورد نظر در جامعه بوده که در آن S مجموعه اعضای جامعه است، یعنی $S = \{Y_{hi} : h = 1, 2, \dots, L; i = 1, 2, \dots, N_h\}$ همچنین فرض کنید $\hat{\theta} = \hat{\theta}(s)$ برآوردگر پارامتر مورد نظر در نمونه باشد که در آن $s = \{y_{hi} : h = 1, 2, \dots, L; i = 1, 2, \dots, n_h\}$ میانگین جامعه را با μ نشان داده و برآوردگر ناریب آن را با $\bar{Y}_{st} = \sum_{h=1}^L W_h \bar{Y}_h$ نمایش می دهیم که در آن $W_h = N_h/N$ و $\bar{Y}_h = \sum_{i=1}^{n_h} Y_{hi}/n_h$ ککران [۳] برآوردگری ناریب برای $Var(\bar{Y}_{st})$ به صورت زیر ارائه

باشد. یک الگوی نمونه گیری طبقه ای با جایگذاری با حجم نمونه n_h در طبقه h م در این جامعه به اجرا در می آوریم. با این فرضیات و با استفاده از ویژگی های نمونه گیری طبقه ای با جایگذاری می توان واریانس میانگین نمونه را به صورت زیر تعریف کرد:

$$Var(\bar{Y}) = \sum_{i=1}^L W_h^2 \frac{\sigma_h^2}{n_h} \quad (۳)$$

که در آن

$$\sigma_h^2 = \frac{1}{n_h} \sum_{i=1}^{n_h} (y_{ih} - \bar{y}_h)^2.$$

از طرفی

$$\sigma_h^2 = \frac{n_h - 1}{n_h} s_h^2. \quad (۴)$$

با جایگذاری رابطه (۴) در رابطه (۳) اثبات قضیه کامل می شود. □

با مقایسه رابطه (۴) با رابطه (۳) می توان دریافت که برآوردگری که با استفاده از بوت استرپ استاندارد حاصل می شود یک برآوردگر اریب برای $Var(\bar{Y})$ است. نتیجه مهمتری که می توان به آن اشاره کرد این است که این برآوردگر برای واریانس سازگار نیست [۹]. البته این مشکل را می توان با در نظر گرفتن $n_h = n_0$ و $f_h = f$ به ازای هر h حل کرد. در این صورت عبارت زیر

$$\frac{n_0(1-f)v_b}{n_0-1}$$

یک برآوردگر نااریب سازگار برای $Var(\bar{Y})$ خواهد بود. روش های دیگری نیز برای حل این مشکل ارائه شده است که در ادامه به طور مختصر به آنها اشاره کرده و به برخی از ویژگی های آنها اشاره می کنیم.

یا تقریب مونت کارلوی آن یعنی:

$$\tilde{v}_B = \frac{1}{B} \sum_{b=1}^B (\hat{\theta}_b^* - \bar{\theta}^*)^2 \quad ; \quad \bar{\theta}^* = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^*.$$

که در آن امید ریاضی باز نمونه گیری می باشد. در معادلات فوق $E_*\bar{\theta}^*$ و $\bar{\theta}^*$ را می توان با $\hat{\theta}$ جایگزین کرد.

نکته ۱ لازم به ذکر است که امید ریاضی بوت استرپ (E_*) با استفاده از اصل جایگزینی^{۱۰} یعنی با جایگزین کردن تابع توزیع تجربی به جای تابع توزیع جامعه حاصل می شود.

نکته ۲ الگوریتم های مختلف باز نمونه گیری فقط در شیوه انتخاب نمونه بوت استرپ باهم تفاوت داشته و بقیه مراحل آنها یکسان هستند.

قضیه ۱ : در ساده ترین حالت وقتی $\hat{\theta} = \bar{y}$ است می توان دید.

$$v_b = \sum_{h=1}^L \frac{W_h^2}{n_h} \left(\frac{n_h - 1}{n_h} \right) s_h^2. \quad (۲)$$

اثبات. الگوی باز نمونه گیری که در الگوریتم قبل بیان شد را می توان به صورت زیر تعبیر کرد: نمونه ای که در اختیار داریم با اجرای یک نمونه گیری طبقه ای بدون جایگذاری از جامعه اصلی که در آن حجم طبقه h م برابر N_h و حجم نمونه حاصل از آن طبقه برابر n_h است، بدست آمده است. حال این نمونه را به عنوان جامعه جدید در نظر می گیریم که جامعه ای است متشکل از L طبقه که حجم جامعه در طبقه h م برابر n_h می

^{۱۰}Plog-in principle

۱.۳ بوت استرپ بدون جایگذاری

ابتدا این روش را در یک نمونه تصادفی ساده بدون جایگذاری بیان کرده سپس آنرا در نمونه گیری طبقه ای تعریف می کنیم. بوت استرپ بدون جایگذاری ابتدا توسط گروس [۵] برای برآورد واریانس معرفی شد. فرض کنید $N = kn$ باشد که k عددی صحیح است، الگوریتم روش BWO به صورت زیر است:

(۱) تولید شبیه جامعه ای توسط k بار تکرار نمونه اصلی.

(۲) اجرای نمونه گیری بدون جایگذاری به حجم n از این جامعه تا حصول نمونه بوت استرپ y_1^*, \dots, y_n^* .

(۳) بدست آوردن $\hat{\theta}^* = \hat{\theta}(y_1^*, \dots, y_n^*)$.

پس از محاسبه پارامتر بوت استرپ، بقیه الگوریتم شبیه به الگوریتم بوت استرپ استاندارد است.

قضیه ۲: وقتی $\hat{\theta} = \bar{y}$ است در روش گروس می توان دید.

$$v_{bwo} = \frac{k(n-1)}{kn-1} \cdot \frac{1-f}{n} s^2. \quad (5)$$

اثبات. فرض کنید x_1, \dots, x_n نمونه اصلی و y_1, \dots, y_{kn} شبه جامعه حاصل از k بار تکرار نمونه باشد. همانطور که اشاره شد نمونه تصادفی بدون جایگذاری از جامعه استخراج می کنیم. لذا براساس نمونه گیری بدون جایگذاری میانگین و واریانس این نمونه عبارتست از:

$$\bar{y} = \frac{1}{kn} \sum_{i=1}^{kn} y_i = \frac{k}{kn} \sum_{i=1}^n x_i = \bar{x}.$$

$$\begin{aligned} Var(\bar{y}) &= \frac{1-f}{n} s^2 \\ &= \frac{1-f}{n} \cdot \frac{1}{kn-1} \sum_{i=1}^{kn} (y_i - \bar{x})^2 \\ &= \frac{1-f}{n} \cdot \frac{k}{kn-1} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \frac{1-f}{n} \cdot \frac{k(n-1)}{kn-1} s^2 \end{aligned}$$

□ که اثبات قضیه کامل می شود.

همانطوری که ملاحظه می شود روش BWO حتی در ساده ترین حالت قادر به تولید برآوردگر ناریب برای واریانس میانگین نمونه نیست. البته در این حالت خاص نیز با استفاده از یک ضریب تصحیح می توان برآوردگر ناریب به دست آورد، اما این امر در الگوهای پیچیده تر نمونه گیری به راحتی انجام نمی گیرد. بیگل و فریدمن [۲] روشی را تحت عنوان بسط روش BWO در نمونه گیری طبقه ای ارائه دادند که بعدها توسط مک کارتی و اسنودن [۸] تکمیل گردید. روش آنها به این صورت است که فرض کنید برای هر طبقه رابطه $N_h = k_h n_h + r_h$ برقرار باشد که در آن $0 \leq r_h \leq n_h - 1$ و k_h و n_h اعداد صحیحی هستند. در این حالت دو شبه جامعه را طوری تشکیل می دهیم که یکی با k_h بار تکرار طبقه h ام و دیگری با $k_h + 1$ بار تکرار آن حاصل می شود. سپس از هر طبقه به طور مستقل نمونه ای بدون جایگذاری به حجم n_h با احتمال p_h از جامعه اول و با احتمال $1 - p_h$ از جامعه دوم استخراج می کنیم. که

$$p_h = \frac{(1-f_h)/(n_h-1) - a_{h2}}{a_{h1} - a_{h2}}.$$

که در آن

$$a_{hi} = \frac{k_h + i - 2}{(k_h + i - 1)n_h - 1}; \quad i = 1, 2.$$

کردن استفاده می شود. به این صورت که تعریف می کنیم $n'_{h1} = n_h - 1$, $k_{h2} = [k_h]$, $k_{h1} = \lceil k_h \rceil$ و $n'_{h1} = n_h$ که در آن $[.]$ و $\lceil . \rceil$ به ترتیب جزء صحیح و جزء صحیح به اضافه یک هستند. حال در رویه گفته شده با احتمال p_h از زوج (n'_{h1}, k_{h1}) و با احتمال $1 - p_h$ از زوج (n'_{h2}, k_{h2}) استفاده می کنیم. که

$$p_h = \frac{(1 - f_h)/(n_h - n_h) - a_{h2}}{a_{h1} - a_{h2}},$$

که در آن

$$a_{hi} = \frac{k_{hi}(1 - n'_{hi}/n_h.k_{hi})}{n'_{hi}(n_h.k_{hi} - 1)}, \quad i = 1, 2.$$

در این حالت نیز سیترا ثابت کرده است که برآوردگر نارایب برای واریانس میانگین نمونه حاصل خواهد شد.

۲.۳ باز مقیاس گذاری بوت استرپ

بسیاری از آماره ها قابل نوشتن به صورت $\hat{\theta} = g(\bar{y})$ یعنی تابعی از میانگین نمونه هستند. ضرایب رگرسیون، نسبت میانگین ها و ضریب همبستگی از این نوع اند. رانو و وو [۹] برای برآورد واریانس یک برآوردگر مانند $\hat{\theta}$ روشی با الگوریتم زیر در نظر گرفتند:

(۱) استخراج نمونه ای با جایگذاری از نمونه اصلی.

(۲) باز مقیاس گذاری هر عضو باز نمونه گیری شده.

(۳) محاسبه برآوردگر اصلی توسط بردار مقادیر باز مقیاس گذاری شده.

فاکتورهای باز مقیاس گذاری طوری انتخاب می شوند که واریانس باز نمونه گیری با واریانس معمولی در حالت

مک کارتی و اسنودن [۸] ثابت کردند که در این حالت واریانس برآوردگر $\bar{y}_{st}^* = \sum_{h=1}^L W_h \bar{y}_h^*$ عبارتست از

$$v_{bwo}(\bar{y}_{st}^*) = \sum_{h=1}^L W_h^2 \frac{1 - f_h}{n_h} s_h^2.$$

که همان $Var(\bar{y}_{st})$ در نمونه گیری طبقه ای معمولی است. روش مک کارتی و اسنودن به راحتی قابل اجراء است اما متأسفانه تحت شرایطی ممکن است p_h عددی منفی شود. مثالی از این حالت را مک کارتی و اسنودن [۸] برای $h = 1$ ارائه دادند. برای رفع این مشکل سیترا [۱۰] نیز روشی تحت عنوان بوت استرپ بدون جایگذاری در نمونه گیری طبقه ای ارائه داد که اساس کار او نیز مشابه روش مک کارتی و اسنودن [۸] تولید دو شبه جامعه است اما به شیوه ای متفاوت. روش سیترا اینگونه است که برای طبقه h ام تعریف می کنیم:

$$n'_h = n_h - (1 - f_h),$$

$$k_h = \frac{N_h}{n_h} \left(1 - \frac{1 - f_h}{n_h}\right).$$

فرض کنید n'_h و k_h اعداد صحیحی باشند. نمونه های هر طبقه را مستقلاً k_h بار تکرار می کنیم تا شبه جامعه ای حاصل شود. سپس نمونه هایی به حجم n'_h بدون جایگذاری از هر طبقه استخراج می کنیم. پس از استخراج نمونه بوت استرپ بقیه مراحل مشابه قبل انجام می شود. تحت این شرایط سیترا ثابت کرده است که برآوردگر واریانس بوت استرپ برابر با واریانس نمونه در حالت معمولی خواهد بود. یعنی این روش برآوردگری نارایب برای واریانس میانگین نمونه تولید می کند.

به وضوح می توان دید که ممکن است n'_h و k_h اعداد صحیحی نباشند. در این صورت از رویه تصادفی

۳.۳ بوت استرپ باجایگذاری

این روش که دارای ساختار ساده ای است توسط مک کارتی و اسنودن [۸] معرفی شد. در الگوی نمونه گیری طبقه ای روش بوت استرپ باجایگذاری (*BWR*) به این صورت است که نمونه ای تصادفی باجایگذاری در هر طبقه به حجم نمونه اصلی در آن طبقه به طور مستقل انتخاب می شود. پس استخراج نمونه بوت استرپ بقیه مراحل مشابه قبل انجام می شود. از مزایای این روش می توان به سازگار بودن آن با ساختار بوت استرپ که نمونه ها به صورت باجایگذاری انتخاب می شوند اشاره کرد. این روش حالت خاصی از روش *M.M*^{۱۱} است که بعدها توسط سیتز [۱۰] معرفی شد و در ادامه مقاله به آن پرداخته می شود. بنابراین تمام مواردی که برای روش *M.M* گفته می شود در این الگو نیز صادق است.

۴ روش *M.M*

طراحی الگوی های بازنمونه گیری بهتر است به گونه ای باشد که با الگوی نمونه گیری اصلی داده ها متناسب باشند. به عبارتی اگر نمونه اصلی بدون جایگذاری انتخاب شده است بهتر است بازنمونه گیری هم طوری طراحی شود که نمونه ها به صورت بدون جایگذاری انتخاب شوند و به طور مشابه اگر نمونه اصلی باجایگذاری باشد منطقی است که بازنمونه گیری هم باجایگذاری باشد. این مسئله در مورد روش *BWO* صادق است چراکه در این روش از شبه جامعه تولید شده با کسر نمونه گیری برابر با کسر نمونه گیری داده های

برآوردگرهای خطی برابر باشد.

در نمونه گیری طبقه ای اگر $\{y_{hi}, i = 1, \dots, n_h\}$ مقادیر نمونه اصلی در طبقه h ام و $\{y_{hi}^*, i = 1, \dots, n_h^*\}$ نمونه ای تصادفی ساده از آن باشد، عامل بازمقیاس گذاری به صورت زیر تعریف می شود:

$$C_h = \sqrt{n_h^*} (n_h - 1)^{-1/2} (1 - f_h)^{1/2}.$$

با استفاده از این عامل مقادیر بازنمونه گیری به صورت زیر بازمقیاس گذاری می شوند:

$$\tilde{y}_h^* = \bar{y}_h + C_h (\bar{y}_h^* - \bar{y}_h)$$

که در آن \bar{y}_h و \bar{y}_h^* به ترتیب میانگین نمونه اصلی و بازنمونه گیری در طبقه h ام هستند. پس از این مرحله میانگین بازمقیاس گذاری، یعنی $\tilde{y}^* = \sum_{h=1}^L W_h \tilde{y}_h^*$ محاسبه شده و سپس برآوردگر اصلی با استفاده از این میانگین بدست می آید ($\hat{\theta}^* = g(\tilde{y}^*)$).

رئو و وو [۹] ثابت کردند که با این الگوریتم برآوردگر ناریب برای توابع خطی از میانگین نمونه قابل دستیابی است. همچنین آنها الگوریتم های مختلفی را برای کلاس بزرگی از مدل های نمونه گیری از جمله خوشه ای دو مرحله ای و نمونه گیری با احتمالات نابرابر ارائه دادند. البته این روش معایبی دارد که برخی از مهمترین آنها عبارتند از:

(۱) محاسبه فاکتور بازنمونه گیری در الگوهای پیچیده نمونه گیری می تواند مشکل ساز باشد.

(۲) این روش حافظ دامنه تغییرات پارامتر نیست.

ثابت می شود که روش $M.M$ بر خلاف روش باز مقیاس گذاری حافظ دامنه تغییرات پارامتر است. (مراجعه شود به [۱۰، ۱۹]).

در ادامه قضیه ای در مورد برآوردگر بوت استرپ حاصل از روش فوق بیان و اثبات می کنیم. ابتدا امید ریاضی و واریانس بوت استرپ که در این قضیه مورد استفاده قرار می گیرند را تعریف می کنیم.

تعریف ۱: فرض کنید $s^* = \{y_{ij}^*; i = 1, \dots, n, j = 1, \dots, B\}$ یک نمونه بوت استرپ از نمونه ای به حجم n با B بار تکرار باشد. همچنین فرض کنید $\theta_j^* = \theta(y_j^*)$ تابعی از نمونه بوت استرپ z ام باشد. امید ریاضی و واریانس بوت استرپ پارامتر θ^* را به ترتیب با نمادهای $E_*(\theta^*)$ و $Var_*(\theta^*)$ نشان داده و با استفاده از تعریف تابع توزیع تجربی به صورت زیر تعریف می شوند.

$$E_*(\theta^*) = \frac{1}{B} \sum_{b=1}^B \theta_b^*,$$

$$var_*(\theta^*) = \frac{1}{B} \sum_{b=1}^B (\theta_b^* - E_*(\theta^*))^2.$$

قضیه ۳: در روش $M.M$ در الگوی نمونه گیری تصادفی ساده داریم:

الف) برآوردگر بوت استرپ یک برآوردگر ناریب برای

$$E_*(\bar{y}^*) = \bar{y} \text{ یعنی نمونه است یعنی } E_*(\bar{y}^*) = \bar{y}.$$

$$b) \text{ } Var_*(\bar{y}^*) = Var(\bar{y})$$

اثبات. براساس تعریف ۱ می توان نوشت

$$E_*(\bar{y}^*) = E_* \left(\frac{1}{n'k} \sum_{i=1}^k \sum_{j=1}^{n'} y_{ij}^* \right)$$

اصلی نمونه هایی به حجم نمونه اصلی و به روش بدون جایگذاری انتخاب می شوند. این مطلب را می توان به عنوان یک مزیت برای این روش دانست. اما در برخی مدل های باز نمونه گیری مانند روش های BWR و یا باز مقیاس گذاری این گونه نیست چرا که در این الگوها نمونه گیری با جایگذاری انجام می شود. اما می توان گفت که این روش ها با ساختار اصلی بوت استرپ سازگار است. زیرا همان طور که می دانیم در روش بوت استرپ نمونه گیری به صورت با جایگذاری از نمونه اصلی صورت می گیرد. حال به نظر می رسد مدلی که این دو ویژگی را توأمأ دارا باشد مدل مناسبی خواهد بود. روشی که در ادامه معرفی می شود دارای این ویژگی است.

۱.۴ روش $M.M$ در الگوی تصادفی ساده

فرض کنید از جامعه ای به حجم N نمونه ای تصادفی ساده بدون جایگذاری به حجم n استخراج شده باشد بطوریکه داشته باشیم $N = k.n$ یا به عبارتی $k = f^{-1}$. تعریف می کنیم $n' = f.n$ ، با فرض صحیح بودن n و k روش $M.M$ برای این الگو عبارتست از:

۱) نمونه گیری تصادفی ساده بدون جایگذاری به

حجم n' از نمونه اصلی تا حصول $y_1^*, y_2^*, \dots, y_{n'}^*$.

۲) ثبت و سپس جایگذاری این زیر نمونه در نمونه

اصلی.

۳) k بار تکرار مستقل مراحل ۱ و ۲ تا حصول نمونه

بوت استرپ n تایی $\underline{y}^* = (y_1^*, \dots, y_n^*)$ و سپس

$$\hat{\theta}^* = \theta(\underline{y}^*)$$

n انتخاب می کنیم. در این حالت الگوریتم زیر در نظر گرفته می شود:

(۱) انتخاب $1 < n' < n$ و بازنمونه گیری به حجم n' تا حصول $y_1^*, \dots, y_{n'}^*$.

(۲) $k = [n(1 - f^*)] / [n'(1 - f)]$ بار (با شرط صحیح بودن k) تکرار مرحله ۱ به طور مستقل و جایگزین کردن نمونه های به حجم n' در هر مرحله تا حصول $y_1^*, \dots, y_{n^*}^*$ که در آن $n^* = n' \cdot k$.

به الگوی فوق روش $M.M$ اصلاح شده گفته می شود.

قضیه ۴: تحت شرایط فوق (الف و ب) نیز نتایج قضیه ۳ در الگوی تصادفی ساده معتبر است. اثبات. اثبات این قضیه به راحتی با نوشتن تعریف امید ریاضی و واریانس بوت استرپ نتیجه می شود. □

۲.۴ روش $M.M$ در الگوی طبقه ای

از آنجایی که الگوی نمونه گیری طبقه ای حالت کلی تر الگوی تصادفی ساده است یا به عبارتی نمونه تصادفی ساده همان نمونه گیری طبقه ای با $L = 1$ است، لذا تمام شرایطی که در قسمت قبل ارائه شد قابل تعمیم به حالت طبقه ای می باشد. فرض کنید در هر طبقه رابطه $N_h = n_h k_h$ (به عبارتی $k_h = f_h^{-1}$) برقرار باشد. تعریف می کنیم $n'_h = f_h n_h$ و فرض می کنیم k_h و n_h به ازای هر h اعداد صحیح باشند. (در غیر این صورت از رویه تصادفی کردن استفاده می کنیم که در قسمت قبل اشاره شد.) آنگاه روش $M.M$ اصلاح شده را برای الگوی طبقه ای به صورت زیر می توان تعریف کرد:

از آنجایی که استخراج k زیر نمونه به طور مستقل صورت می گیرد داریم:

$$E_*(\bar{y}^*) = \frac{1}{k} \sum_{i=1}^k E_* \left(\frac{1}{n'} \sum_{j=1}^{n'} y_{ij}^* \right) = \frac{1}{k} \sum_{i=1}^k \bar{y} = \bar{y}.$$

برای اثبات قسمت دوم قضیه بطور مشابه می توان نوشت:

$$\begin{aligned} Var_*(\bar{y}^*) &= Var_* \left(\frac{1}{n'k} \sum_{i=1}^k \sum_{j=1}^{n'} y_{ij}^* \right) \\ &= \frac{1}{k^2} \sum_{i=1}^k Var_* \left(\frac{1}{n'} \sum_{j=1}^{n'} y_{ij}^* \right) \\ &= \frac{1 - f^*}{k \cdot n'} s^2. \end{aligned}$$

چون کسر بازنمونه گیری با کسر نمونه گیری نمونه اصلی برابر است داریم:

$$Var_*(\bar{y}^*) = \frac{1 - f}{n} s^2 = Var(\bar{y}).$$

□ این اثبات قضیه را کامل می کند. شرطی که در استفاده از این رویه وجود دارد و شاید آنرا بتوان یک محدودیت دانست، لزوم صحیح بودن n' و برقراری شرط $n' = fn \geq 1$ است. البته راه حل هایی برای آن نیز ارائه شده است.

(الف) فرض کنید $n' = fn \geq 1$ بوده ولی عدد صحیح نباشد، آنگاه می توان با عمل تصادفی کردن حجم بازنمونه گیری را به این صورت انتخاب کرد که با احتمال p ، n' را $[fn]$ در نظر گرفته و با احتمال $1 - p$ آنرا $[fn] + 1$ در نظر می گیریم. که در آن p طوری انتخاب می شود که نتیجه برآورد واریانس \bar{y} برابر با برآوردگر نارایب آن در حالت معمولی باشد.

(ب) حال اگر $0 < fn < 1$ باشد، n' را عددی بین ۱ و

(۱) نمونه گیری تصادفی ساده بدون جایگذاری به

حجم $1 < n'_h < n_h$ از نمونه اصلی در هر طبقه تا حصول $y_{h n'_h}^*, \dots, y_{h 2}^*, y_{h 1}^*$

(۲) ثبت و سپس جایگذاری این زیر نمونه در نمونه اصلی.

(۳) $k_h = [n_h(1 - f_h^*)]/[n'_h(1 - f_h)]$ بار تکرار

مستقل مراحل ۱ و ۲ تا حصول نمونه

بوت استرپ $y_{h n_h}^*, \dots, y_{h 2}^*, y_{h 1}^*$ و در نتیجه

$$s^* = \{y_{hi}^* : h = 1, 2, \dots, L; i = 1, 2, \dots, n_h\}$$

(۴) محاسبه $\hat{\theta}^* = \hat{\theta}(s^*)$.

باید توجه داشت که در قسمت نمونه گیری بدون جایگذاری کسر باز نمونه گیری برابر است با کسر نمونه گیری اصلی، به همین دلیل این مرحله را k_h بار تکرار می کنیم تا بردار باز نمونه گیری به حجم n_h حاصل شود.

قضیه ۵: در روش $M.M$ اصلاح شده در الگوی طبقه

ای وقتی $\theta = \bar{Y}$ و $\hat{\theta} = \bar{y}_{st}$ باشند، برای برآوردگر

بوت استرپ آن یعنی $\bar{y}_{st}^* = \sum_{h=1}^L W_h \bar{y}_h^*$ که در آن

$$\bar{y}_h^* = \sum_{i=1}^{k_h} \sum_{j=1}^{n'_h} y_{hij}^* / n_h^*$$

الف) \bar{y}_{st}^* برآوردگر ناریب \bar{y}_{st} است، یعنی

$$E_*(\bar{y}_{st}^*) = \bar{y}_{st}$$

ب) واریانس بوت استرپ برآوردگر بوت استرپ \bar{y}_{st}^*

برابر است با واریانس برآوردگر معمولی \bar{y}_{st} . یعنی

$$Var_*(\bar{y}_{st}^*) = Var(\bar{y}_{st})$$

اثبات. می توان نوشت

$$\begin{aligned} E_*(\bar{y}_h^*) &= E_* \left[\frac{1}{n_h^*} \sum_{i=1}^{k_h} \sum_{j=1}^{n'_h} y_{hij}^* \right] \\ &= \frac{1}{k_h} \sum_{i=1}^{k_h} E_* \left[\frac{1}{n'_h} \sum_{j=1}^{n'_h} y_{hij}^* \right] = \frac{1}{k_h} \sum_{i=1}^{k_h} \bar{y}_h = \bar{y}_h. \end{aligned}$$

پس

$$E_*(\bar{y}_{st}^*) = E_* \sum_{h=1}^L W_h \bar{y}_h^* = \sum_{h=1}^L W_h \bar{y}_h = \bar{y}_{st}.$$

که اثبات الف کامل می شود. برای اثبات قسمت دوم از

آنجایی که نمونه گیری در طبقات مختلف به طور مستقل

انجام می گیرد، داریم:

$$\begin{aligned} Var_*(\bar{y}_{st}^*) &= \sum_{h=1}^L W_h^2 Var_* \left(\frac{1}{k_h n'_h} \sum_{i=1}^{k_h} \sum_{j=1}^{n'_h} y_{hij}^* \right) \\ &= \sum_{h=1}^L \frac{W_h^2}{h_h^2} \sum_{i=1}^{k_h} \left(\frac{1 - f_h^*}{n'_h} s_h^2 \right) \\ &= \sum_{h=1}^L W_h^2 \frac{1 - f_h}{n_h} s_h^2 \end{aligned}$$

□ که اثبات قضیه کامل می شود.

قضیه فوق برای هر n'_h و k_h برقرار است لذا می توان n'_h

را طوری کرد که کسر باز نمونه گیری $(f_h^* = n'_h/n_h)$

برابر با کسر نمونه گیری اصلی $(f_h = n_h/n_h)$ باشد. در

ادامه مزیت این انتخاب توضیح داده خواهد شد.

۱.۲.۴ انتخاب n'_h برای جور شدن گشتاور سوم

در الگوریتم قبل همانطوری که اشاره شد $E_*(\bar{y}^*) = \bar{y}$ و

$Var_*(\bar{y}^*) = Var(\bar{y})$ است. یعنی دو گشتاور اول توزیع

\bar{y}^* برابر با برآوردگر ناریب دو گشتاور اول \bar{y} هستند. در

این بخش هدف پیدا کردن n'_h است به طوری که گشتاور سوم بوت استرپ یعنی $E_*(\bar{y}^* - \bar{y})^3$ برابر با $\hat{\mu}_3(\bar{y})$ شود.

ابتدا یادآور می شویم که برای $n_h \geq 2$ برآوردگر ناریب $E(\bar{y} - E(\bar{y}))^3$ برابر است با:

$$\hat{\mu}_3(\bar{y}) = \sum_{h=1}^L W_h^3 \frac{(\lambda - f_h)(\lambda - 2f_h)m_{3h}}{(n_h - \lambda)(n_h - 2)},$$

که در آن

$$m_{3h} = \frac{1}{n_h} \sum_{i=1}^{n_h} (y_{hi} - \bar{y}_h)^3.$$

(۷)

گشتاور نمونه ای مرتبه سوم در طبقه h ام است. با فرض صحیح بودن n'_h و k_h ، حجم کل نمونه در هر طبقه عبارتست از

$$n_h^* = k_h n'_h = \frac{n_h(\lambda - f_h^*)}{\lambda - f_h}.$$

که در آن $f_h^* = n'_h/n_h$ کسر بازنمونه گیری است. با فرض $E_*(\varepsilon_h) = 0$ می توان نوشت

و در نتیجه

$$\begin{aligned} E_*(\bar{y}^* - \bar{y})^3 &= E_* \left(\sum_{h=1}^L \varepsilon_h \right)^3 \\ &= \sum_{h=1}^L E_*(\varepsilon_h^3) \\ &= \sum_{h=1}^L W_h^3 E_*(\bar{y}_h^* - \bar{y}_h)^3. \end{aligned} \quad (6)$$

حال برای پیدا کردن $E_*(\bar{y}_h^* - \bar{y}_h)^3$ ، با صرف نظر کردن از زیرنویس h بنا به مستقل بودن طبقات می توان نوشت

$$\bar{y}^* = \frac{1}{k} \sum_{j=1}^k \bar{y}_j^* = \frac{1}{kn'} \sum_{j=1}^k \sum_{l=1}^{n'_j} y_{jl}^*.$$

که در آن زیر نویس های l و j به ترتیب مربوط به قسمت های با جایگذاری و بدون جایگذاری در الگوریتم

هستند. حال داریم:

$$\begin{aligned} E_*(\bar{y}^* - \bar{y})^3 &= \frac{1}{k^3} \left(\sum_{j=1}^k E_*(\bar{y}_j^* - \bar{y})^3 + 3 \sum_{j \neq l} E_*(\bar{y}_j^* - \bar{y}) E_*(\bar{y}_l^* - \bar{y})^2 \right. \\ &\quad \left. + \sum_{j \neq l \neq m} E_*(\bar{y}_j^* - \bar{y}) E_*(\bar{y}_l^* - \bar{y}) E_*(\bar{y}_m^* - \bar{y}) \right) \\ &= \frac{1}{k^3} \sum_{j=1}^k E_*(\bar{y}_j^* - \bar{y})^3 \\ &= \frac{n^3 (\lambda - f^*) (\lambda - 2f^*) m_3}{(kn'^3)(n - \lambda)(n - 2)}. \end{aligned} \quad (۷)$$

رابطه اخیر از این حقیقت که \bar{y}_j^* ها از یک نمونه تصادفی با جایگذاری به حجم n' از نمونه اصلی حاصل می شوند و از نتیجه مقاله سوخاتم و سوخاتم [۱۲] به دست می آید. داریم:

$$\begin{aligned} &\frac{n^3 (\lambda - f^*) (\lambda - 2f^*) m_3}{(kn'^3)(n - \lambda)(n - 2)} \\ &= \frac{(\lambda - f_h) (\lambda - 2f_h) m_{3h}}{(n_h - \lambda)(n_h - 2)} [1 - A_h]. \end{aligned} \quad (۸)$$

که در آن

$$A_h = \frac{(f_h^* - f_h)}{(\lambda - f_h^*)(\lambda - 2f_h)}.$$

با جایگذاری روابط (۷) و (۸) در رابطه (۶) داریم:

$$\begin{aligned} E_*(\bar{y}^* - \bar{y})^3 &= \sum_{h=1}^L W_h^3 \frac{(\lambda - f_h) (\lambda - 2f_h) m_{3h}}{(n_h - \lambda)(n_h - 2)} [1 - A_h]. \end{aligned}$$

A_h در رابطه قبل عامل خطای جور شدن می باشد که اگر کسر بازنمونه گیری با کسر نمونه گیری اصلی برابر باشد $A_h = 0$ خواهد شد. این یعنی اگر $n'_h = f_h = n_h$ باشد گشتاور سوم بازنمونه گیری با گشتاور سوم نمونه یکسان خواهد شد.

۲.۲.۴ تصادفی کردن انتخاب n'_h و k_h در حالت غیر

صحیح

در ابتدا می توان به راحتی نشان داد که در الگوی نمونه گیری طبقه ای اگر $1 \leq n'_h \leq n_h / (2 - f_h)$ باشد آنگاه

$$k_h = [n(1 - f_h^*)] / [n'(1 - f_h)] \geq 1$$

خواهد بود. حال اگر k_h عدد صحیحی نباشد متغیر K_h را به صورت زیر تعریف می کنیم:

$$p_h = P(K_h = k_{h1}) = \left(\frac{1}{k_h} - \frac{1}{k_{h2}} \right) / \left(\frac{1}{k_{h1}} - \frac{1}{k_{h2}} \right)$$

$$P(K_h = k_{h2}) = 1 - p_h$$

که در آن $k_{h1} = [k_h]$ و $k_{h2} = [k_h] + 1$ بوده که در شرط $1 \leq k_{h1} \leq k_h \leq k_{h2} \leq n_h$ صدق می کنند. حال با الگوریتمی که در قسمت قبل ارائه شد و با استفاده از K_h ، رویه را به طور مستقل در هر طبقه انجام می دهیم.

قضیه ۶: با استفاده از روش فوق اگر $\hat{\theta} = \bar{y}$ باشد آنگاه $Var_*(\bar{y}^*) = Var(\bar{y})$ خواهد بود.

اثبات. می دانیم

$$Var_*(\bar{y}^*) = \sum_{h=1}^L W_h^2 Var(\bar{y}_h^*). \quad (9)$$

از آنجایی که نمونه گیری در طبقه ها به طور مستقل انجام می شود، در رابطه فوق $Var(\bar{y}_h^*)$ را برای یک طبقه به دست آورده و قضیه را ثابت می کنیم. لذا با در نظر نگرفتن زیرنویس h داریم:

$$Var_*(\bar{y}^*) = E_{1*} var_{2*}(\bar{y}^*) + Var_{1*} E_{2*}(\bar{y}^*). \quad (10)$$

که در آن E_{1*} و Var_{1*} به ترتیب امید ریاضی و واریانس باز نمونه گیری و E_{2*} و Var_{2*} امید ریاضی و واریانس

باز نمونه گیری به شرط K هستند. از طرفی

$$E_{1*} Var_{2*}(\bar{y}^*)$$

$$= p Var(\bar{y}^* | K = k_1) + (1 - p) Var(\bar{y}^* | K = k_2)$$

$$= p \times \frac{1 - f^*}{k_1 n'} s^2 + (1 - p) \times \frac{1 - f^*}{k_2 n'} s^2$$

$$= \frac{1 - f^*}{n'} \left(\frac{p}{k_1} + \frac{1 - p}{k_2} \right) s^2.$$

با جایگذاری p در عبارت فوق داریم

$$E_{1*} Var_{2*}(\bar{y}^*) = \frac{1 - f^*}{k n'} s^2.$$

از طرفی

$$Var_{1*} E_{2*}(\bar{y}^*) = Var_{1*}(\bar{y}) = 0.$$

که تساوی اخیر از استقلال \bar{y} از k نتیجه می شود. با استفاده از این نتایج و روابط (۹) و (۱۰) اثبات قضیه کامل می شود. □

در روش فوق اگر از فرضیه $n'_h = f_h n_h$ استفاده شود احتمال این وجود دارد که هم n_h و هم k_h غیر صحیح باشند. در این حالت نیز از روش هایی برای تصادفی کردن این دو اختیار است تا تساوی $Var_*(\bar{y}^*) = Var(\bar{y})$ برقرار شود. یک روش تصادفی کردن k به شرط صحیح بودن n' در مرحله قبل اشاره شد. در این قسمت می خواهیم n' را نیز تصادفی در نظر بگیریم. یک روش به این صورت است که فرض کنید $n'_1 = [n']$ ، $n'_2 = [n']$ ، $k_1 = [k]$ و $k_2 = [k]$ باشند، آنگاه جفت (n'_1, k_1) را با احتمال p و (n'_2, k_2) را با احتمال $1 - p$ انتخاب می کنیم که

$$p = \frac{\left(\frac{1-f}{n} - \frac{1-n'_2/n}{k_2 n'_2} \right)}{\left(\frac{1-f}{n} - \frac{1-n'_1/n}{k_1 n'_1} \right) - \left(\frac{1-f}{n} - \frac{1-n'_2/n}{k_2 n'_2} \right)}.$$

دانشگاه امری وقت گیر است، با اجرای یک الگوی نمونه گیری طبقه ای با در نظر گرفتن دانشکده ها به عنوان طبقه، میانگین و انحراف معیار آن را برآورد می کنیم. در این جامعه $N = 210$ ، $L = 10$ و $N_i = \{33, 25, 22, 30, 33, 14, 13, 12, 13, 15\}$ است. با در نظر گرفتن کران خطای نمونه گیری (D) بر اساس رابطه زیر در تعیین حجم نمونه در الگوی طبقه ای، n را مشخص می کنیم.

$$n_0 = \frac{Z_{1-\alpha/2}^2 S^2}{D^2}; \quad n = \frac{n_0}{1 + n_0/N}$$

در رابطه فوق $Z_{1-\alpha/2}^2$ و S^2 به ترتیب چندک $1 - \alpha/2$ توزیع نرمال و واریانس جامعه هستند. با فرض $\alpha = 0.05$ و $T = 35$ و بر اساس نتایج حاصل از بررسی روی داده های این جامعه توسط آقای گنجی [۱] داریم $S^2 = 12840/9$ و لذا خواهیم داشت $n = 70$ و $f = 0.33$. با این فرضیات می توان حجم نمونه ها را با اعمال تخصیص متناسب به صورت زیر در نظر گرفت:

$$n_i = \{11, 8, 7, 10, 11, 5, 4, 4, 4, 6\}$$

بر اساس شرایط فوق با استفاده از چهار روش بازنمونه گیری معرفی شده در این مقاله به برآورد میانگین جامعه و واریانس آن می پردازیم. نتایج حاصله به ازای $B = 1000$ در جدول ۱ ارائه شده است. لازم به ذکر است که نتایج با استفاده از نرم افزار R به دست آمده اند.

با در نظر گرفتن رویه بیان شده در حالتی که $\hat{\theta} = \bar{y}$ با راحتی می توان نشان داد $Var_*(\bar{y}^*) = Var(\bar{y})$. با در نظر نگرفتن زیر نویس h داریم:

$$\begin{aligned} Var_*(\bar{y}^*) &= p Var_*(\bar{y}^* | K = k_1, n' = n'_1) \\ &\quad + (1-p) var_*(\bar{y}^* | K = k_2, n' = n'_2) \\ &= p \times \frac{1 - n'_1/n}{k_1 n'_1} s^2 \\ &\quad + (1-p) \times \frac{1 - n'_2/n}{k_2 n'_2} s^2. \end{aligned}$$

با جایگذاری p در عبارت فوق نتیجه لازم حاصل می شود.

۵ مقایسه عددی

در این قسمت دقت روش های مختلف بازنمونه گیری ارائه شده در این مقاله را برای یک مثال واقعی مقایسه می شوند. برای انجام این کار تولیدات علمی اساتید دانشگاه فردوسی مشهد در نظر گرفته شده اند. جامعه مورد مطالعه جامعه ای است شامل ۱۰ طبقه (دانشکده) و هدف تحقیق برآورد میانگین تولیدات علمی اساتید و تعیین میزان دقت این برآورد با استفاده از محک انحراف معیار است. تولیدات علمی اساتید بر مبنای فعالیت های علمی آنان مانند تألیف کتاب، تألیف مقاله، ترجمه کتاب و مقاله و نظیر آنها محاسبه می شود. از آنجایی که محاسبه تولیدات علمی تمام اساتید

جدول ۱. جدول ارائه شده برای مقایسه روشهای بازنمونه گیری.

| پارامتر | \bar{y}_{rs}^* | \bar{y}_{wr}^* | \bar{y}_{mm}^* | \bar{y}_{bwo}^* | \bar{y}_{st} |
|---------|------------------|------------------|------------------|-------------------|----------------|
| میانگین | ۲۲۵/۸۳ | ۲۲۵/۱۰ | ۲۲۵/۰۱ | ۲۲۵/۱۳ | ۲۲۵/۱۸ |
| واریانس | ۱۴۷/۴۶ | ۱۰۲/۴۹ | ۱۱۰/۰۳ | ۱۰۶/۲۵ | ۱۰۸/۹۸ |

در مورد واریانس میانگین نمونه هم براساس داده های جدول به نظر می رسد که روش های $M.M$ و BWO از دقت بیشتری برخوردار باشند.

۶ نتیجه گیری

در این مقاله روش های مختلف بازنمونه گیری در الگوی طبقه ای ارائه و مقایسه شده اند. با بررسی های انجام شده می توان گفت که روش $M.M$ کاراتر از دیگر روش های ارائه شده می باشد. این روش برای تمام مدل های نمونه گیری (طبقه ای؛ خوشه ای و نمونه گیری با احتمالات نابرابر) به راحتی قابل اجرا بوده و برای کلاس بزرگی از برآوردگرها از جمله واریانس توابع غیر خطی از میانگین و صدک ها معتبر است [۱۱].

با مشاهده این جدول می توان دید که هر چهار روش بازنمونه گیری تقریباً برآورد یکسانی برای میانگین جامعه دارند و این برآورد ها بسیار به میانگین نمونه (\bar{y}_{est}) نزدیک هستند. همان طوری که قبلاً اشاره شد روش های بازنمونه گیری برآوردی برای آماره در نمونه اصلی هستند. به عنوان مثال در این جامعه میانگین برابر ۲۱۷/۵۵ است (براساس تحقیق آقای گنجی) و برآوردگر میانگین جامعه یعنی میانگین نمونه برابر ۲۲۵/۱۸. از جدول فوق پیداست که برآوردگرهای بازنمونه گیری به میانگین نمونه ای نزدیکتر هستند تا به میانگین جامعه. یعنی آنها برآوردگرهای برای میانگین نمونه هستند. حال هرچه برآوردگر نمونه ای برای پارامتر جامعه دقیقتر باشد بازنمونه گیری نیز مطمئن تر خواهند بود.

مراجع

- [۱] گنجی، م. (۱۳۸۱)، مقایسه تولیدات علمی اساتید دانشکده های دانشگاه فردوسی مشهد، رساله کارشناسی ارشد، دانشگاه فردوسی مشهد.
- [2] Bickel, P.J. and Freedman, D.A. (1984), Asymptotic normality and the Bootstrap in stratified sampling., *Ann.Statist.* **12**, 470-482.
- [3] Cochran, W.G. (1977), Sampling Techniques., *John Wiley, New York*.
- [4] Efron, B. (1979), Bootstrap methods: another look ot the jackknife., *Ann.Statist.* **7**, 1-26.
- [5] Gross, S. (1980), Median Estimation in Sample Survey., *Proceedings of Section on Survey Research Methods, American Statistical Association.* , 181-184.

- [6] Krewski, D., and Rao, J. N. K. (1981), Inference From Stratified Samples: Properties of the linearization, Jackknife and Balanced Repeated Replication Methods., *The Annals of Statistics* **9**, 1010-1019.
- [7] Loh, W.Y. (1987), Calibration Confidence Coefficients., *Journal of the American Statistical Association* **82**, 155-162.
- [8] McCarthy, P.J., and Snowden, C.B. (1985), The Bootstrap and finite population sampling., *Vital and Health Statistics* **2**, 1-23.
- [9] Rao, J.N.K., and Wu, C.F.J. (1988), Resampling interference with complex survey data., *Journal of the American Statistical Association* **83**, 231-241.
- [10] Sitter, R.R. (1992a), Comparing three Bootstrap methods for survey data., *The Canadian Journal of Statistics* **20**, 135-154.
- [11] Sitter, R.R. (1992b), Resampling procedures for complex survey data., *Journal of the American Statistical Association* **87**, 135-154.
- [12] Sukhatme, P.V. and Sukhatme, B.W. (1970), Sampling theory of survey with applications., *Iowa State University Press, Ames, IA*
- [13] Wu, C.F.J. (1985), Variance Estimation for the Combined Ratio and Combined Regression Estimators., *Journal of the Royal Statistical Society, Ser.B* **47**, 147-154.