

مقایسه‌ی برآوردهای آنتروپی طرح‌های نمونه‌گیری

فهمیه مسیحی بیدگلی

گروه آمار، دانشکده ریاضی، دانشگاه صنعتی اصفهان

چکیده

یکی از روش‌های مقایسه‌ی طرح‌های نمونه‌گیری محاسبه‌ی میزان آنتروپی آن‌ها است. آنتروپی یک طرح نمونه‌گیری اندازه‌ای از تصادفی بودن و گستردگی آن را نشان می‌دهد. محاسبه‌ی آنتروپی از طریق تعریف، به‌علت زیاد بودن نمونه‌های ممکن کاری بسیار وقت‌گیر و گاهی غیر عملی است. هدف این مقاله، یافتن برآوردهای مناسب آنتروپی برای طرح‌های نمونه‌گیری است. استفاده از یک برآوردها، به طرح و موقعیت آن بستگی دارد. بعضی از برآوردها تنها برای طرح‌هایی که دارای تابع احتمال مشخص هستند، مناسب‌اند و بعضی از آن‌ها را حتی در حالتی که تابع احتمال طرح نمونه‌گیری مشخص نباشد، می‌توان مورد استفاده قرار داد. مقایسه‌ی این برآوردها با استفاده از چند مثال شبیه‌سازی صورت گرفته است و نتایج نشان می‌دهد، در صورت مشخص نبودن تابع احتمال، استفاده از برآوردها آنتروپی خام با تصحیح میزان آریبی مناسب‌تر خواهد بود.

واژه‌های کلیدی: آنتروپی، طرح نمونه‌گیری، نمونه‌گیری پواسون شرطی تعدیل‌شده، نمونه‌گیری پارتو.

۱ مقدمه

تا قبل از آن تنها در علم فیزیک استفاده می‌شد، در سایر علوم جایگاه خود را پیدا کرد. بدین علت آنتروپی به‌طور قراردادی با آنتروپی شانون نمایش داده شد. آنتروپی یک طرح نمونه‌گیری میزان گستردگی و تصادفی بودن یک طرح را نشان می‌دهد. هاجک [۲] برای اولین بار آنتروپی طرح‌های نمونه‌گیری را بررسی کرد و نشان داد در کلاسی از طرح‌های نمونه‌گیری با احتمال

مفهوم آنتروپی در علوم که با پدیده‌های تصادفی در ارتباط هستند، کاربرد دارد و در رشته‌های گوناگون علمی معانی متفاوتی می‌یابد. از جمله معانی که برای آنتروپی به‌کار می‌رود، آشفتگی، بی‌نظمی، عدم قطعیت و میزان تصادفی بودن یک پیشامد است. شانون [۷] معیار آنتروپی را با علم آمار و احتمال پیوند داد و آنتروپی که

به عنوان میزان عدم قطعیتی که برای مقدار X وجود دارد تفسیر کرد. آنتروپی با تعریف فوق آنتروپی شانون نامیده می شود. محققین زیادی سعی در توسعه آنتروپی داشته اند. از جمله، رنی [۴] و تی سالیس [۹] آنتروپی های دیگری را معرفی کرده اند. در ادامه به بعضی از ویژگی های آنتروپی اشاره می شود.

۱- آنتروپی تابعی از احتمال های p_1, \dots, p_n است و به مقادیر x_i ($i = 1, 2, \dots, n$) بستگی ندارد.

۲- آنتروپی تابعی پیوسته است و هر تغییر بسیار اندک در مقدارهای احتمال، باعث تغییر بسیار کوچک در میزان آنتروپی می شود.

۳- تابع آنتروپی نسبت به هر یک از مؤلفه های خود متقارن است و اندازه آنتروپی با تغییر در ترتیب برآمدهای x_i تغییر نخواهد کرد.

۴- اگر همه برآمدهای x_i هم شانس باشند، اندازه آنتروپی بیشینه می شود و در این حالت آنتروپی با افزایش تعداد برآمدها، افزایش می یابد.

۵- آنتروپی توأم دو سیستم مستقل برابر با جمع آنتروپی هر یک از آن ها خواهد بود:

$$H(X, Y) = H(X) + H(Y).$$

در صورتی که دو سیستم مستقل نباشند، آنتروپی توأم عبارت است از

$$H(X, Y) = H(X) + H(Y|X),$$

که در آن $H(Y|X)$ متوسط عدم قطعیتی است که بعد از دریافت اطلاعات X ، برای Y باقی می ماند و بر اساس رابطه زیر به دست می آید.

$$H(Y|X) = \sum_{i=1}^n H(Y|X = x_i) Pr(X = x_i),$$

برابر و با اندازه نمونه ثابت، طرح نمونه گیری تصادفی ساده بدون جایگذاری دارای حداکثر میزان آنتروپی است. او همچنین نشان داد در کلاس طرح های نمونه گیری با احتمال نابرابر و با بردار احتمال شمول مرتبه اول یکسان و با اندازه نمونه ثابت، طرح نمونه گیری پواسون شرطی تعدیل شده^۱ حداکثر آنتروپی را دارد. گرافسترم [۱] نیز آنتروپی طرح های نمونه گیری با احتمال نابرابر را مورد بررسی قرار داد و به این نتیجه رسید که چندین طرح در داشتن حداکثر آنتروپی به یکدیگر نزدیک هستند. همچنین او چندین برآوردگر آنتروپی را در موقعیت های مختلف معرفی کرد که بعضی از این برآوردها در این مقاله با یکدیگر مقایسه می شوند. ابتدا در بخش ۲ آنتروپی و کاربرد آن در طرح نمونه گیری بیان و سپس در بخش ۳ برآوردهای مختلف آنتروپی طرح نمونه گیری ارائه می شود. در بخش ۴ با بررسی دو جامعه و شبیه سازی نمونه های طرح های مختلف، برآوردهای آنتروپی مقایسه شده و در انتها نتایج ارائه می شود.

۲ آنتروپی

اگر X متغیر تصادفی باشد که یکی از مقادیر x_1, \dots, x_n را با احتمال های p_1, \dots, p_n انتخاب کند، $-\log p_i$ نشان دهنده میزان تعجب حاصل از آن است که مقدار x_i را اختیار کند، در نتیجه امید ریاضی میزان تعجب از اطلاع در مورد متغیر تصادفی X برابر است با

$$H(X) = E[-\log Pr(X)] = -\sum_{i=1}^n p_i \log p_i.$$

کمیت $H(X)$ در نظریه اطلاعات، به عنوان آنتروپی متغیر تصادفی X شناخته می شود. می توان $H(X)$ را

^۱ Adjusted Conditional Poisson

که در آن

$$H(Y|X = x_i) = - \sum_j Pr(y_j|x_i) \log Pr(y_j|x_i).$$

که در آن x به عنوان نمونه‌ی تصادفی در نظر گرفته شده است. در صورتی که جرم احتمال روی مجموعه نمونه‌های ممکن به خوبی توزیع شود، آنتروپی طرح نمونه‌گیری مقدار بالایی خواهد داشت و هنگامی که تنها تعداد کمی از نمونه‌ها احتمال زیادی داشته باشند، مقدار آنتروپی کم خواهد بود. از دیگر عواملی که در مقدار آنتروپی طرح نمونه‌گیری تأثیرگذار است، اندازه‌ی تکیه‌گاه طرح است. در صورت زیاد بودن اندازه‌ی تکیه‌گاه، آنتروپی نیز مقدار بالایی خواهد داشت.

۱.۲ آنتروپی طرح نمونه‌گیری

همان‌طور که اشاره شد، آنتروپی یک طرح نمونه‌گیری اندازه‌ای از میزان تصادفی بودن طرح را نشان می‌دهد. هنگامی که طرح نمونه‌گیری آنتروپی بالایی دارد، یک میزان زیاد از عدم قطعیت یا میزان بالایی از تعجب در نمونه‌ی انتخابی وجود دارد. به عبارت دیگر، در این حالت میزان بالایی از تصادفی بودن وجود دارد. اگر اطلاعات اضافی در مورد جامعه و ویژگی‌های آن موجود باشد، این اطلاعات می‌تواند برای انتخاب طرح به کار رود. در صورت عدم دسترسی به این اطلاعات، طرح با آنتروپی بالا به علت تصادفی بودن زیاد، انتخاب مناسبی برای انجام نمونه‌گیری خواهد بود.

۳ برآورد آنتروپی در طرح نمونه‌گیری

در صورتی که اندازه‌ی نمونه‌ی n و اندازه جامعه N باشد حداکثر $\binom{N}{n}$ نمونه‌ی ممکن برای یک طرح نمونه‌گیری بدون جایگذاری با اندازه‌ی نمونه‌ی n وجود دارد. بنابراین در صورت بزرگ بودن اندازه‌ی جامعه، به علت وجود تعداد نمونه‌های بسیار زیاد محاسبه‌ی آنتروپی از طریق تعریف بسیار وقت‌گیر است. گاه نیز امکان دارد تابع احتمال طرح نمونه‌گیری مشخص نباشد، از این رو به برآورد آنتروپی از طریق شبیه‌سازی روی آورده می‌شود. در صورتی که تابع احتمال یک طرح نمونه‌گیری را نتوان مشخص کرد با استفاده از الگوریتم نمونه‌گیری، نمونه‌های طرح شبیه‌سازی تولید و تابع احتمال برآورد می‌شود. یک الگوریتم نمونه‌گیری روشی است که برای انتخاب یک نمونه صرف نظر از تابع احتمال آن به کار می‌رود. برای برآورد تابع احتمال بدون هیچ اطلاعاتی، زمان زیادی صرف می‌شود.

اگر m تعداد نمونه‌های شبیه‌سازی شده و m_i برای i

یک طرح نمونه‌گیری یک توزیع گسسته روی مجموعه‌ای از نمونه‌های ممکن \mathbf{x}_k ($k \in Q$) است. به تکیه‌گاه طرح گفته می‌شود و مجموعه‌ای از تمام نمونه‌های ممکن برای یک طرح نمونه‌گیری است. نمونه‌ی \mathbf{x}_k می‌تواند به عنوان برداری از نشانگر عضویت مطرح شود. بنابراین $\mathbf{x}_k \in \{0, 1\}^N$ که N اندازه‌ی جامعه، ۱ نشان‌دهنده‌ی عضویت واحد جامعه در نمونه و ۰ عدم عضویت آن واحد است. احتمال به دست آمدن نمونه‌ی \mathbf{x}_k با $p(\mathbf{x}_k)$ نمایش داده می‌شود و به عنوان طرح نمونه‌گیری شناخته می‌شود که دارای شرط $\sum_{k \in Q} p(\mathbf{x}_k) = 1$ است [۸]. آنتروپی طرح نمونه‌گیری به صورت زیر تعریف می‌شود.

$$H = - \sum_{k \in Q} p(\mathbf{x}_k) \log p(\mathbf{x}_k) = -E_p[\log(p(\mathbf{x}))],$$

است. در این وضعیت از برآوردگر \hat{H}_{PF} به شرح زیر استفاده می‌شود.

اگر \mathbf{x}_k برای $k = 1, 2, \dots, m$ نمونه‌ی شبیه‌سازی شده‌ی k -ام و $p(\mathbf{x}_k)$ احتمال این نمونه بر اساس تابع طرح مشخص باشد، آنگاه برآوردگر

$$\hat{H}_{PF} = -\frac{1}{m} \sum_{k=1}^m \log(p(\mathbf{x}_k)),$$

برای برآورد آنتروپی طرح نمونه‌گیری به‌کار می‌رود. این برآوردگر ناریب است و برای طرح‌هایی که دارای تابع احتمال مشخص هستند، استفاده می‌شود. اگر بتوان تابع احتمال p را توسط p_0 تقریب زد، یک برآوردگر مناسب دیگر از آنتروپی ایجاد می‌شود.

$$\hat{H}_{APF\backslash} = -\frac{1}{m} \sum_{k=1}^m \log(p_0(\mathbf{x}_k)).$$

برآوردگر $\hat{H}_{APF\backslash}$ اریب است. می‌توان با محاسبه‌ی $E_p\left[\frac{p(\mathbf{x})}{p_0(\mathbf{x})}\right]$ میزان اریبی این برآوردگر را کم کرد. برای این منظور ابتدا باید $p(\mathbf{x})$ برآورد شود. و برآوردگر آنتروپی به‌صورت زیر حاصل می‌شود:

$$\begin{aligned} \hat{H}_{APF\backslash} &= -\frac{1}{m} \sum_{k=1}^m \log(p_0(\mathbf{x}_k)) \\ &- \frac{1}{2} \left(\frac{1}{m} \sum_{k=1}^m \frac{\hat{p}(\mathbf{x}_k)}{p_0(\mathbf{x}_k)} - 1 \right). \end{aligned}$$

۴ شبیه‌سازی

برای مقایسه‌ی برآوردهای آنتروپی از مثال‌های شبیه‌سازی استفاده می‌کنیم. شبیه‌سازی با استفاده از دو جامعه‌ی معروف ترات-باندسون-میستر و سامفورد-هاجک انجام شده است. در این شبیه‌سازی از نمونه‌گیری‌های پارتو [۵، ۶] و پواسون شرطی

$1, \dots, m$ نشان‌دهنده‌ی تعداد تکرارهای i -امین نمونه‌ی شبیه‌سازی شده باشد، برآوردگر حداکثر درست‌نمایی احتمال نمونه‌های ممکن برابر $\frac{m_i}{m}$ بوده و یک برآوردگر خام^۲ آنتروپی به‌صورت زیر خواهد بود.

$$\hat{H}_{naive} = -\frac{1}{m} \sum_{i=1}^m \log\left(\frac{m_i}{m}\right).$$

می‌توان نشان داد که برآوردگر فوق دارای اریبی منفی است. هر چه تعداد نمونه‌های شبیه‌سازی بیشتر باشد، اریبی این برآوردگر کمتر و دقت برآورد آن بیشتر خواهد بود. این برآوردگر به برآوردگر حداکثر درست‌نمایی نیز معروف است.

اگر تعداد نمونه‌های متمایز، q و برآوردگر حداکثر درست‌نمایی احتمال نمونه‌های ممکن $\hat{p}(\mathbf{x}_i) = \frac{m_i}{m}$ باشد، شکل دیگر نمایش برآوردگر خام آنتروپی عبارت است از

$$\hat{H}_{naive} = -\sum_{i=1}^q \hat{p}(\mathbf{x}_i) \log(\hat{p}(\mathbf{x}_i)).$$

محققین مختلفی در جهت رفع اریبی این برآوردگر تلاش کرده‌اند. تصحیح اریبی^۳ زیر برای برآوردگر خام، توسط میلر [۳] پیشنهاد شد.

$$\hat{H}_{BC} = -\frac{1}{m} \sum_{i=1}^m \log\left(\frac{m_i}{m}\right) + \frac{q}{2m}.$$

اگر تابع احتمال معلوم و برای هر نمونه‌ی ممکن قابل محاسبه باشد، می‌توان با انتخاب m به اندازه‌ی کافی بزرگ و یک برآوردگر مناسب، برآورد دقیق آنتروپی طرح نمونه‌گیری را به‌دست آورد. هنگامی که اندازه‌ی جامعه و نمونه بزرگ باشد، محاسبه‌ی آنتروپی از طریق تعریف به علت بررسی تعداد زیاد نمونه‌های ممکن، بسیار وقت‌گیر

^۲Naive estimator

^۳Bias Correction

مقدار واقعی بسیار کمتر برآورد می‌کند. همچنین فاصله‌ی اطمینان برای این برآوردگر در حالتی که تعداد نمونه ۵۰۰ است مقدار واقعی آنتروپی را در بر نرفته است. اما با تصحیح میزان اریبی، این برآوردگر بسیار بهتر عمل کرده و حتی با تعداد نمونه‌ی کم برآوردی نزدیک به مقدار واقعی خود ارائه می‌دهد. در واقع در برآوردگر \hat{H}_{BC} همانند برآوردگر \hat{H}_{PF} با افزایش تعداد نمونه برآوردها دقیق‌تر و به مقدار واقعی نزدیک‌تر می‌شود.

همان‌طور که پیش از این توضیح داده شد برآوردهای $\hat{H}_{APF\gamma}$ و $\hat{H}_{APF\lambda}$ در حالتی به‌کار می‌رود که توزیع طرح نمونه‌گیری مشخص باشد و یا با طرح دیگری بتوان تقریب زد. اما در بسیاری از روش‌های نمونه‌گیری امکان تعیین توزیع طرح وجود ندارد که در این صورت برآوردگر \hat{H}_{BC} پیشنهاد می‌شود. این برآوردگر بدون نیاز به توزیع طرح نمونه‌گیری با دقتی همانند سه برآوردگر دیگر به‌خوبی عمل می‌کند. اما برای تعداد نمونه‌ی زیاد دارای سرعت برآورد بسیار کمتری است.

مثال ۲. جامعه‌ی سامفورد-هاجک با اندازه‌ی جامعه $N = ۱۰$ ، اندازه‌ی نمونه‌ی $n = ۵$ و بردار احتمال شمول زیر در نظر بگیرید:

$$\pi = (0/2, 0/25, 0/35, 0/4, 0/50, 0/5, 0/55, 0/65, 0/7, 0/9).$$

برای تولید نمونه از نمونه‌گیری پواسون شرطی تعدیل شده استفاده می‌شود. تابع احتمال برای نمونه‌گیری مذکور به‌صورت زیر است.

$$p_{Acp}(\mathbf{x}) = C_{Acp} \prod_{i=1}^{10} (Ap)_i^{x_i} (1 - (Ap)_i)^{1-x_i}, \quad |\mathbf{x}| = 5,$$

که $(Ap)_i$ پارامترهای اصلاح شده و C_{Acp} ثابت نرمال‌ساز

تعدیل شده [۲] استفاده شده است. نتایج به‌دست آمده در هر دو مورد مشابه و مستقل از جامعه و روش نمونه‌گیری است.

مثال ۱. جامعه‌ی ترات-باندسون-میستر با اندازه‌ی جامعه $N = ۶$ ، اندازه‌ی نمونه‌ی $n = ۳$ و بردار احتمال شمول (احتمال آن که واحد k -ام در نمونه‌ی تصادفی قرار گیرد) زیر را در نظر بگیرید:

$$\pi = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{2}{3}, \frac{2}{3}, \frac{2}{3} \right).$$

در این جامعه با استفاده از یک طرح نمونه‌گیری بدون جایگذاری حداکثر $\binom{6}{3} = ۲۰$ نمونه ایجاد می‌شود. برای تولید نمونه از نمونه‌گیری پارتو استفاده می‌شود. تابع احتمال برای نمونه‌گیری مذکور به‌صورت زیر است.

$$p_{par}(\mathbf{x}) = \prod_{i=1}^6 \lambda_i^{x_i} (1 - \lambda_i)^{1-x_i} \times \sum_{k=1}^6 c_k x_k, \quad |\mathbf{x}| = 3.$$

مقدار دقیق آنتروپی با استفاده از بررسی تمام نمونه‌های ممکن $۲/۶۹۵۵۵۸$ است. شبیه‌سازی با تعداد نمونه‌های مختلف برای ۵ برآوردگر مختلف ارائه شده در بخش قبل صورت گرفته است. در دو برآوردگر $\hat{H}_{APF\gamma}$ و $\hat{H}_{APF\lambda}$ برای تقریب p_0 ، تابع احتمال طرح CP (پواسون شرطی) به‌کار می‌رود. به تعداد $n = ۲۰$ بار شبیه‌سازی تکرار و میانگین، انحراف معیار و فاصله اطمینان برآوردهای مختلف در هر بار محاسبه شده و نتایج در جدول ۱ ارائه شده است. با بررسی جدول مشاهده می‌شود، در مورد هر برآوردگر با افزایش تعداد نمونه‌ی شبیه‌سازی برآورد آنتروپی دقیق‌تر و به مقدار واقعی خود نزدیک‌تر می‌شوند. تمام برآوردها به‌جز برآوردگر خام آنتروپی به خوبی و تقریباً مشابه به هم عمل می‌کنند. برآوردگر خام دارای اریبی منفی است و در صورت کم بودن تعداد نمونه از

است. در این جامعه با استفاده از نمونه‌گیری فوق مقدار دقیق آنتروپی $4/726990$ است. در این مثال نیز شبیه‌سازی با تعداد نمونه‌های مختلف برای ۵ برآوردگر صورت گرفته است. در دو برآوردگر \hat{H}_{APF1} و \hat{H}_{APF2} ، برای تقریب p_0 ، تابع احتمال طرح $APar$ (پارتوی تعدیل شده) به کار می‌رود. به تعداد $n = 10$ بار شبیه‌سازی تکرار شده است. نتایج در جدول ۲ مشاهده می‌شود. در این مثال روش نمونه‌گیری و نوع جامعه تغییر داده شد. اما باز نتایج مشابه مثال قبل حاصل می‌شود. یعنی در مورد هر برآوردگر با افزایش تعداد نمونه برآوردها به مقدار واقعی نزدیک‌تر می‌شوند و خطای برآورد مقدار کمتری را نشان می‌دهند. برآوردگر خام دارای اریبی منفی است و با تصحیح میزان اریبی دارای دقت بسیار بهتری می‌شود. برای استفاده از برآوردگر خام باید تعداد نمونه‌های شبیه‌سازی شده بسیار زیاد باشد تا این برآوردگر برآوردی نزدیک به مقدار واقعی آنتروپی داشته باشد. برآوردگرهای آنتروپی در حالتی که طرح نمونه‌گیری دارای تابع احتمال مشخص یا تقریبی می‌باشند تقریباً مشابه هم عمل کرده و استفاده از برآوردگر خام با تصحیح میزان اریبی دارای خطای برآورد کمتر است و به کارگیری آن در صورت مشخص نبودن تابع احتمال طرح نمونه‌گیری پیشنهاد می‌شود.

جدول ۱: مقایسه‌ی برآوردگرهای مختلف آنترویی در جامعه‌ی TBM و با استفاده از نمونه‌گیری پارتو

فاصله اطمینان	انحراف معیار	میانگین	تعداد نمونه	برآوردگر
(۶۹۲۳۳/۲، ۶۶۹۸۰/۲)	۰۲۵۷/۰	۶۸۱۰۷/۲	۵۰۰	\hat{H}_{naive}
(۶۹۷۵۷/۲، ۶۷۶۸۹/۲)	۰۲۳۶/۰	۶۸۷۲۳/۲	۱۰۰۰	
(۶۹۵۰۵/۲، ۶۸۴۷۹/۲)	۰۰۸۳/۰	۶۸۹۹۲/۲	۵۰۰۰	
(۷۰۹۷۷/۲، ۶۹۰۰۰/۲)	۰۲۲۵/۰	۶۹۹۸۹/۲	۵۰۰	\hat{H}_{BC}
(۷۰۶۷۵/۲، ۶۹۲۱۰/۲)	۰۱۶۷/۰	۶۹۹۴۲/۲	۱۰۰۰	
(۶۹۷۳۴/۲، ۶۹۱۶۹/۲)	۰۰۶۴/۰	۶۹۴۵۱/۲	۵۰۰۰	
(۷۱۰۸۲/۲، ۶۸۸۷۲/۲)	۰۲۵۲/۰	۶۹۹۷۷/۲	۵۰۰	\hat{H}_{PF}
(۷۰۴۸۲/۲، ۶۸۹۵۴/۲)	۰۱۷۴/۰	۶۹۷۱۷/۲	۱۰۰۰	
(۶۹۸۱۱/۲، ۶۹۲۷۳/۲)	۰۰۶۱/۰	۶۹۵۴۲/۲	۵۰۰۰	
(۷۱۱۳۱/۲، ۶۸۸۵۱/۲)	۰۲۶۱/۰	۶۹۹۹۱/۲	۵۰۰	\hat{H}_{APF1}
(۷۰۸۴۵/۲، ۶۸۹۱۳/۲)	۰۱۵۸/۰	۶۹۸۷۹/۲	۱۰۰۰	
(۷۰۰۴۰/۲، ۶۹۲۸۰/۲)	۰۰۶۵/۰	۶۹۶۶۰/۲	۵۰۰۰	
(۷۱۰۶۵/۲، ۶۸۹۰۱/۲)	۰۲۳۵/۰	۶۹۹۸۳/۲	۵۰۰	\hat{H}_{APF2}
(۷۰۵۷۰/۲، ۶۸۹۳۶/۲)	۰۱۶۹/۰	۶۹۷۵۳/۲	۱۰۰۰	
(۶۹۹۰۹/۲، ۶۹۲۷۹/۲)	۰۰۵۹/۰	۶۹۵۹۴/۲	۵۰۰۰	

جدول ۲: مقایسه‌ی برآوردگرهای مختلف آنترویی در جامعه‌ی سامفورد-هاجک و با استفاده از نمونه‌گیری پواسون شرطی تعدیل شده

فاصله اطمینان	انحراف معیار	میانگین	تعداد نمونه	برآوردگر
(۶۲۳۵۹/۴، ۵۹۵۷۲/۴)	۰۲۲۵/۰	۶۰۹۶۵۴/۴	۱۰۰۰	\hat{H}_{naive}
(۶۹۵۴۲/۴، ۶۸۱۹۳/۴)	۰۱۰۹/۰	۶۸۸۶۷/۴	۵۰۰۰	
(۷۱۸۹۴/۴، ۷۰۵۶۷/۴)	۰۱۰۷/۰	۷۱۲۳۱/۴	۱۰۰۰۰	
(۷۴۹۹۸/۴، ۷۳۴۹۹/۴)	۰۱۲۱/۰	۷۴۲۴۸/۴	۱۰۰۰	\hat{H}_{BC}
(۷۳۲۳۵/۴، ۷۲۲۰۸/۴)	۰۰۸۳/۰	۷۲۷۲۲/۴	۵۰۰۰	
(۷۲۷۸۳/۴، ۷۱۹۱۲/۴)	۰۰۷۰/۰	۷۲۳۴۷/۴	۱۰۰۰۰	
(۷۵۹۴۳/۴، ۷۳۵۱۶/۴)	۰۱۹۶/۰	۷۴۷۲۹/۴	۱۰۰۰	\hat{H}_{PF}
(۷۳۹۳۲/۴، ۷۲۴۵۸/۴)	۰۱۱۹/۰	۷۳۱۹۵/۴	۵۰۰۰	
(۷۲۶۹۶/۴، ۷۱۹۷۷/۴)	۰۰۵۸/۰	۷۲۳۳۶/۴	۱۰۰۰۰	
(۷۳۵۴۹/۴، ۷۵۸۲۹/۴)	۰/۰ ۲۱۱	۷۴۶۸۹/۴	۱۰۰۰	\hat{H}_{APF1}
(۷۲۳۹۸/۴، ۷۳۵۹۴/۴)	۰۱۰۸/۰	۷۲۹۹۶/۴	۵۰۰۰	
(۷۲۳۲۱/۴، ۷۲۹۰۵/۴)	۰۰۶۵/۰	۷۲۶۱۳/۴	۱۰۰۰۰	
(۷۵۵۳۱/۴، ۷۳۵۳۷/۴)	۰۲۰۵/۰	۷۴۵۳۴/۴	۱۰۰۰	\hat{H}_{APF2}
(۷۳۶۵۲/۴، ۷۲۳۹۶/۴)	۰۱۲۰/۰	۷۳۰۲۴/۴	۵۰۰۰	
(۷۲۹۴۲/۴، ۷۲۱۹۸/۴)	۰۰۶۲/۰	۷۲۵۷۰/۴	۱۰۰۰۰	

York.

- [3] Miller, G. (1955). Note on the bias of information estimates. In: H. Quastler (Ed.), *Information Theory in Psychology II-B*, Free Press, Glencoe, IL, 95-100.
- [4] Renyi, A. (1961). On measures of entropy and information. *Proc, Berkeley Symposium, Statist, Probability*, 1, 547-561.
- [5] Rosen, B. (1997). Asymptotic theory for order sampling. *J. Statist. Plann. Inference*, 62, 135-158.
- [6] Rosen, B. (1997). On sampling with probability proportional to size. *J. Statist. Plann. Inference*, 62, 159-191.
- [7] Shannon, C.E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27, 379-423, 623-656.
- [8] Tille, Y. (2006). *Sampling Algorithms*. Springer series in statistics, Springer science, Business media, Inc., New York.
- [9] Tsallis, C. (1988). Possible generalizations of Boltzmann-Gibbs statistics. *Journal of Statistical Physics*, 52, 479-487.

۵ بحث و نتیجه‌گیری

در صورتی که اندازه‌ی نمونه n و اندازه‌ی جامعه N باشد حداکثر $\binom{N}{n}$ نمونه‌ی ممکن برای یک طرح نمونه‌گیری بدون جایگذاری با اندازه‌ی نمونه‌ی n وجود دارد. بنابراین در صورت بزرگ بودن اندازه‌ی جامعه و نمونه، به‌علت وجود تعداد نمونه‌های بسیار زیاد محاسبه‌ی آنتروپی از طریق تعریف بسیار وقت‌گیر است. در این صورت با استفاده از شبیه‌سازی می‌توان آنتروپی طرح نمونه‌گیری را برآورد کرد. در این مقاله برآوردهای مختلف آنتروپی با یکدیگر مقایسه شد. هر چه تعداد نمونه‌های شبیه‌سازی شده بیشتر باشند برآوردهای انجام شده توسط برآوردهای مختلف به مقدار واقعی نزدیک‌تر خواهند بود. به‌خصوص دربارهی برآوردهای خام باید تعداد نمونه‌های شبیه‌سازی شده بسیار زیاد باشد تا برآورد به مقدار واقعی نزدیک شود. این برآوردهای دارای اریبی منفی بوده و کمتر از مقدار واقعی برآورد می‌کند. با تصحیح میزان اریبی برآوردهای خام، برآورد به مقدار واقعی بسیار نزدیک‌تر می‌شود و در صورت نامشخص بودن تابع توزیع طرح نمونه‌گیری استفاده از این برآوردهای مناسب خواهد بود. در حالتی که توزیع طرح نمونه‌گیری مشخص است برآوردهای \hat{H}_{PF} که دارای سرعت برآورد بالاتری نسبت به برآوردهای \hat{H}_{BC} است، پیشنهاد می‌شود.

مراجع

- [1] Grafstrom, A. (2010). Entropy of unequal probability sampling designs. *Statist. Methodol.*, 7, 84-97.
- [2] Hajek, J. (1981). *Sampling From a Finite Population*. Marcel Dekker, New