

مجله علوم آماری، بهار و تابستان ۱۳۹۰

جلد ۵، شماره ۱، ص ۴۱-۶۰

برآورد بیزی و اعتبارسنجی متقابل پهنای باند برآوردگر هسته‌ای تابع چگالی برای داده‌های در طول-اریب

مسعود عجمی بختیاروند^۱، وحید فکور^۱، سارا جمهوری^۲

^۱گروه آمار، دانشگاه فردوسی مشهد

^۲گروه آمار، دانشگاه بیرجند

تاریخ دریافت: ۱۳۹۰/۲/۱۰ تاریخ آخرین بازنگری: ۱۳۹۰/۵/۱۶

چکیده: چنانچه در نمونه‌گیری، داده‌ها با احتمالی متناسب با اندازه انتخاب شوند، داده‌های حاصل را در طول-اریب نامند. برآورد ناپارامتری تابع چگالی با استفاده از داده‌های در طول-اریب، مشکل‌تر از سایر حالات است. یکی از برآوردگرهای معروف در این زمینه توسط جونز (۱۹۹۱) معرفی شده است. در این مقاله ابتدا پارامتر پهنای باند این برآوردگر با رهیافت بیزی برآورد می‌شود. سپس سازگاری قوی آن با به کار بردن پهنای باند برآورد شده به روش بیزی اثبات می‌شود. در انتها با مطالعه شبیه سازی به مقایسه عملکرد روش بیزی و اعتبارسنجی متقابل در برآورد پهنای باند پرداخته می‌شود.

واژه‌های کلیدی: برآوردگر هسته‌ای چگالی، پهنای باند، داده‌های در طول-اریب، اعتبارسنجی متقابل کمترین توان‌های دوم، هسته گاوسی وارون.

آدرس الکترونیک مسئول مقاله: مسعود عجمی بختیاروند، ajami.masoud@yahoo.com
کد موضوع‌بندی ریاضی (۲۰۰۰): ۶۲G۰۷

پدیده در طول-اریب^۱، اولین بار توسط ویکسل (۱۹۲۵) در آناتومی مطرح شد که وی نام آن را مسأله گلبولی نهاد. او در هنگام دیدن گلبول‌ها در میکروسکوپ متوجه شد که فقط گلبول‌هایی قابل مشاهده هستند که اندازه بزرگی آن‌ها، از حد معینی بیشتر باشد و گلبول‌های کوچکتر قابل دیدن در میکروسکوپ نیستند. بعدها این موضوع توسط مک فادن (۱۹۶۲)، بلومنتال (۱۹۶۷) و کاکس (۱۹۶۹) به مفهوم آماری مورد بررسی قرار گرفت. کاکس (۱۹۶۹) در نمونه گیری یک نوع محصول صنعتی، متوجه شد که الیاف‌های با طول بلندتر، با احتمال بیشتری وارد نمونه می‌شوند. این موضوع نوعی اریبی را به نتایج تحمیل کرد که به "در طول-اریبی" معروف شد. در حالت کلی، اگر در نمونه گیری، عناصری از جامعه که اندازه، طول یا عمر بیشتری نسبت به بقیه اعضا جامعه دارند، وارد نمونه شوند یا شانس ورود به نمونه آن‌ها بیشتر باشد یا به عبارت دیگر عناصر جامعه با احتمالی متناسب با اندازه طول‌شان وارد نمونه شوند، نمونه دچار نوعی اریبی به نام در طول-اریبی می‌شود.

تعریف ۱: فرض کنید F یک تابع توزیع تجمعی مطلقاً پیوسته با تابع چگالی f باشد. متغیر تصادفی Y را در طول-اریب گویند هرگاه تابع توزیع آن به صورت

$$G(y) = \int_0^y \frac{t}{\mu} dF(t), \quad y \geq 0, \quad (1)$$

باشد، که در آن $\mu = \int_0^{\infty} tf(t)dt < \infty$.

با توجه به رابطه (۱) چگالی متغیر تصادفی Y به صورت زیر است

$$g(y) = \frac{yf(y)}{\mu}, \quad y \geq 0. \quad (2)$$

در واقع رابطه (۲) بیانگر آن است که احتمال مشاهده هر مقدار متغیر تصادفی Y متناسب با طول آن است. معمولاً تابع چگالی f را چگالی نااریب و g را چگالی اریب گویند. در مسأله برآورد تابع چگالی f با داده‌های در طول-اریب، افراد زیادی تحقیق نموده اند که از آن جمله می‌توان به باتاچاریا و همکاران (۱۹۸۹)،

^۱ Length-biased

جونز (۱۹۹۱)، گیامون و همکاران (۱۹۹۸)، افروموویچ (۲۰۰۴) و همچنین چوبی و همکاران (۲۰۱۰) اشاره کرد. کریستوبال و آکالا (۲۰۰۱) مرور خوبی بر تحقیقات انجام شده در زمینه داده‌های در طول-اریب انجام داده‌اند. فرض کنید Y_n, \dots, Y_1 نمونه‌ای تصادفی از توزیع G باشد. آنگاه با توجه به رابطه (۴)، یک برآوردگر طبیعی برای $f(x)$ ، به صورت

$$\tilde{f}_n(x) = \frac{\hat{\mu}_n \hat{g}_n(x)}{x}, \quad x > 0, \quad (3)$$

است، که در آن

$$\hat{g}_n(x) = \frac{1}{nh} \sum_{j=1}^n K\left(\frac{x - Y_j}{h}\right) \quad (4)$$

برآوردگر هسته‌ای^۲ تابع چگالی g است که توسط روزن بلات (۱۹۵۶) معرفی شد. در (۴) تابع نامنفی $K(\cdot)$ را تابع هسته می‌نامند که در شرایط

$$\begin{aligned} i) & \int_{-\infty}^{\infty} K(x) dx = 1, \\ ii) & \int_{-\infty}^{\infty} xK(x) dx = 0, \\ iii) & \sigma^2 = \int_{-\infty}^{\infty} x^2 K(x) dx \neq 0 \end{aligned}$$

صدق می‌کند و h مقداری مثبت است که پهنای باند^۳ نامیده می‌شود. همچنین $\hat{\mu}_n$ یک برآوردگر مناسب پارامتر μ است. با توجه به اینکه

$$E\left(\frac{1}{\bar{Y}}\right) = \int_0^{\infty} \frac{1}{y} g(y) dy = \frac{1}{\mu},$$

می‌توان μ را با روش گشتاوری به صورت

$$\hat{\mu}_n = \frac{n}{\sum_{i=1}^n Y_i^{-1}}$$

برآورد کنیم. برآوردگر $\tilde{f}_n(x)$ توسط باتاچاریا و همکاران (۱۹۸۸) معرفی و خواص حدی آن از قبیل سازگاری ضعیف و نرمال بودن مجانبی مورد بررسی قرار گرفت.

^۲ Kernel estimator

^۳ Bandwidth

۴۴ برآورد بیزی و اعتبار سنجی متقابل پهنای باند برآوردگر هسته‌ای تابع چگالی

مشکل اساسی این برآوردگر این است که وقتی $x \rightarrow \infty$ آنگاه $\tilde{f}_n(x) \rightarrow \infty$ برای رفع این مشکل جونز (۱۹۹۱) برآوردگر

$$\hat{f}_n(x) = (nh)^{-1} \hat{\mu}_n \sum_{j=1}^n K\left(\frac{x - Y_j}{h}\right) Y_j^{-1}, \quad (5)$$

را معرفی کرد، که در آن $K(\cdot)$ یک هسته متقارن است. جونز (۱۹۹۱)، ضمن بررسی خواص بهینه این برآوردگر نشان داد این برآوردگر یک تابع چگالی احتمال نیز است. بعلاوه در نزدیکی صفر مشکل برآوردگر باتاچاریا و همکاران (۱۹۸۸) را ندارد و رفتار بهتری از خود نشان می‌دهد. وی با استفاده از شاخص $MSE(\hat{f}_n(x)) = \int (\hat{f}_n(x) - f(x))^2 dx$ ، h بهینه را به دست آورد. شکل کلی این برآوردگر برای توابع هسته دلخواه (متقارن یا نامتقارن)، به صورت زیر است

$$\hat{f}_h(x) = n^{-1} \hat{\mu}_n \sum_{j=1}^n Y_j^{-1} K(x; Y_j, h). \quad (6)$$

در برآورد تابع چگالی به روش هسته دو عامل مهم، تابع هسته و پارامتر پهنای باند نقش مهمی را ایفا می‌کنند. هر چند لازم به یادآوری است که انتخاب پارامتر پهنای باند بسیار مهمتر از انتخاب تابع هسته است. منابع متعددی در این زمینه وجود دارند که از آن جمله می‌توان به روزن بلات (۱۹۵۶)، پارزن (۱۹۶۲)، سیلورمن (۱۹۸۵) و واند و جونز (۱۹۹۵) اشاره نمود. بطور کلی توابع هسته را می‌توان به دو دسته‌های متقارن و نامتقارن تقسیم کرد، که از جمله هسته‌های متقارن می‌توان به هسته‌های باکس کار^۴، گاوسی و اپانچ‌نیکوف^۵ اشاره کرد. همچنین هسته‌های گاما، گاوسی وارون^۶ و عکس گاوسی وارون^۷ از جمله هسته‌های نامتقارن هستند.

معمولاً هسته‌های نامتقارن در برآورد تابع چگالی با تکیه‌گاه محدود مورد استفاده قرار می‌گیرند. زیرا هسته‌های متقارن در برآورد این نوع توابع چگالی اغلب خوب

^۴ Boxcar kernel

^۵ Epanechnikov kernel

^۶ Inverse Gussian kernel

^۷ Reciprocal inverse Gussian kernel

مسعود عجمی بختیاروند، وحید فکور، سارا جمهوری ۴۵

عمل نمی‌کنند و مشکل اریبی مرزی^۸ را بوجود می‌آورند (کولاسکرا و پاچت، ۲۰۰۶).

چون h میزان همواری برآوردگر را کنترل می‌کند به آن پارامتر همواری نیز گفته می‌شود. با انتخاب پهنای بانده کوچک، برآورد کم هموار و با پهنای بانده بزرگ‌تر، برآورد بیش همواری برای تابع چگالی به دست می‌آید. در بحث برآورد تابع چگالی، معمولاً پارامتر پهنای بانده مجهول است و باید برآورد شود. چندین روش برای برآورد h وجود دارد، که از آن جمله می‌توان به روش‌های نزدیک‌ترین همسایگی^۹، اعتبارسنجی متقابل کمترین توان‌های دوم^{۱۰} و اعتبارسنجی متقابل اریب^{۱۱} اشاره کرد.

در تمامی این روش‌های h به صورت یک پارامتر فراموضعی^{۱۲} در نظر گرفته می‌شود. گن‌گوبادهایا و چه‌آنگ (۲۰۰۲) با در نظر گرفتن رهیافت بیزی به برآورد موضعی h پرداختند. برای داده‌های سانسور شده کولاسکرا و پاچت (۲۰۰۶) با در نظر گرفتن تابع زیان توان دوم خطا و انتخاب پیشین مناسب برای پهنای بانده، برآوردگر $h(x)$ را به صورت میانگین توزیع پسین به دست آوردند.

در این مقاله، در بخش ۲ به بررسی انتخاب پهنای بانده به روش بیزی در برآوردگر جونز پرداخته می‌شود. اثبات قضایای این بخش در پیوست ارائه می‌شوند. در بخش ۳ با مطالعه شبیه سازی دو برآورد بیزی و اعتبارسنجی متقابل کمترین توان‌های دوم برای پهنای بانده در برآوردگر جونز مقایسه می‌شوند. بحث و نتیجه‌گیری در بخش ۴ ارائه می‌شود.

۲ برآورد بیزی پهنای بانده

در این قسمت، با فرض اینکه پهنای نوار خود یک متغیر تصادفی است و داده‌ها در طول-اریب هستند، از برآوردگر جونز برای تخمین زدن تابع چگالی f استفاده

^۸ Boundary bias

^۹ Nearest neighbour

^{۱۰} Least squares cross validation

^{۱۱} Biased cross validation

^{۱۲} Global

۴۶ برآورد بیزی و اعتبار سنجی متقابل پهنای باند برآوردگر هسته‌ای تابع چگالی

می‌شود. این مسأله را با انتخاب دو هسته گاوسی وارون و گاوسی که به ترتیب توابعی نامتقارن و متقارن هستند، بررسی نموده و h به روش بیزی برآورد می‌شود. فرض کنید $\pi(h)$ تابع چگالی پیشین برای h باشد. با در نظر گرفتن نمونه تصادفی Y_1, \dots, Y_n از توزیع G ، چگالی پسین h در نقطه x به صورت زیر است

$$\hat{\pi}(h|Y_1, \dots, Y_n, x) = \frac{\hat{f}_h(x)\pi(h)}{\int \hat{f}_h(x)\pi(h)dh} \quad (V)$$

در قضیه زیر، پهنای باند برآوردگر جونز به صورت موضعی و بر اساس تابع هسته گاوسی وارون و چگالی پیشین گاما وارونه محاسبه می‌شود. شکل کلی تابع هسته گاوسی وارون به صورت زیر است

$$K(x; y, h) = \frac{1}{\sqrt{\pi}hx^{\alpha}} e^{-\frac{1}{h}(x-y)^2/xy^2}, \quad x, h > 0, -\infty < y < \infty. \quad (A)$$

قضیه ۱: فرض کنید تابع هسته، گاوسی وارون و توزیع پیشین h ، گامای وارون با پارامترهای α و β به صورت

$$\pi(h) = \frac{1}{\beta^{\alpha}\Gamma(\alpha)h^{\alpha+1}} e^{-\frac{1}{\beta h}}, \quad \alpha > 0, \beta > 0. \quad (9)$$

باشد. با در نظر گرفتن تابع زیان توان دوم خطا، برآورد بیزی h که میانگین توزیع پسین است در نقطه x برابر است با

$$h_n^* = h_n(x) = E[(h|Y_1, \dots, Y_n, x)] = \frac{\sum_{j=1}^n Y_j^{-1} (\beta_j^*)^{(\alpha^*-1)}}{(\alpha^* - 1) \sum_{j=1}^n Y_j^{-1} (\beta_j^*)^{(\alpha^*)}}. \quad (10)$$

برهان: با جایگذاری تابع هسته (A) در رابطه (9) و با توجه به (V) توزیع پسین h به صورت

$$\hat{\pi}(h|Y_1, \dots, Y_n, x) = \frac{\sum_{j=1}^n Y_j^{-1} \cdot \frac{1}{\sqrt{\pi}hx^{\alpha}} e^{-\frac{(x-Y_j)^2}{hxy^2}} \cdot \frac{1}{\beta^{\alpha}\Gamma(\alpha)h^{\alpha+1}} e^{-\frac{1}{\beta h}}}{\int \sum_{j=1}^n Y_j^{-1} \cdot \frac{1}{\sqrt{\pi}hx^{\alpha}} e^{-\frac{(x-Y_j)^2}{hxy^2}} \cdot \frac{1}{\beta^{\alpha}\Gamma(\alpha)h^{\alpha+1}} e^{-\frac{1}{\beta h}} dh}$$

$$= \frac{\sum_{j=1}^n (Y_j^{-1} / h^{\alpha^*+1} e^{\frac{-1}{\beta_j^* h}})}{\sum_{j=1}^n Y_j^{-1} \int (e^{\frac{-1}{\beta_j^* h}} / h^{\alpha^*+1}) dh}$$

$$= \frac{\sum_{j=1}^n (Y_j^{-1} / h^{\alpha^*+1} e^{\frac{-1}{\beta_j^* h}})}{\Gamma(\alpha^*) \sum_{j=1}^n Y_j^{-1} (\beta_j^*)^{\alpha^*}},$$

به دست می آید، که در آن

$$\beta_j^* = \left[\frac{1}{\beta} + \frac{(x - Y_j)^2}{2xY_j^2} \right]^{-1}$$

$$\alpha^* = \alpha + \frac{1}{\beta}, \quad \alpha > \frac{1}{\beta}.$$

بر آورد بیزی h میانگین توزیع پسین به صورت زیر است.

$$h_n^* = \int h \cdot \frac{\sum_{j=1}^n (Y_j^{-1} / h^{\alpha^*+1}) e^{-1/(h\beta_j^*)}}{\Gamma(\alpha^*) \sum_{j=1}^n Y_j^{-1} (\beta_j^*)^{\alpha^*}} dh$$

$$= \frac{\sum_{j=1}^n Y_j^{-1}}{\Gamma(\alpha^*) \sum_{j=1}^n Y_j^{-1} (\beta_j^*)^{\alpha^*}} \int \frac{e^{-1/h\beta_j^*}}{h^{\alpha^*}} dh$$

$$= \frac{\sum_{j=1}^n Y_j^{-1} (\beta_j^*)^{(\alpha^*-1)}}{(\alpha^* - 1) \sum_{j=1}^n Y_j^{-1} (\beta_j^*)^{\alpha^*}}.$$

قضیه ۲: فرض کنید Y_1, \dots, Y_n نمونه‌ای تصادفی از توزیع G ، تابع چگالی f کراندار و $T = \inf\{t : 1 - F(t) > 0\}$ باشد. اگر به ازای هر $0 < x < T$ ، وقتی

$$\lim_{n \rightarrow \infty} h_n^* = 0$$

آنگاه

$$|\hat{f}_{h_n^*}(x) - f(x)| \xrightarrow{a.s.} 0.$$

برهان: با در نظر گرفتن G_n به عنوان تابع توزیع تجربی Y_i ها می توان تابع توزیع F را به صورت زیر برآورد کرد

$$F_n(t) = \hat{\mu}_n \int_0^t y^{-1} dG_n(y).$$

بنابراین

$$\begin{aligned} \hat{f}_{h_n^*}(x) &= n^{-1} \hat{\mu}_n \sum_{j=1}^n Y_j^{-1} K(x; Y_j, h_n^*) \\ &= \hat{\mu}_n \int_0^\infty u^{-1} K(x; u; h_n^*) dG_n(u) \\ &= \int_0^\infty K(x; u; h_n^*) dF_n(u). \end{aligned} \quad (11)$$

در نظر بگیرید

$$f_{h_n^*}(x) = \int_0^\infty K(x; u; h_n^*) dF(u).$$

بنابراین

$$\begin{aligned} |\hat{f}_{h_n^*}(x) - f(x)| &\leq |\hat{f}_{h_n^*}(x) - f_{h_n^*}(x)| + |f_{h_n^*}(x) - f(x)| \\ &=: J_1 + J_2. \end{aligned} \quad (12)$$

حال داریم

$$\begin{aligned} J_1 &= \left| \int_0^\infty K(x; u; h_n^*) dF_n(u) - \int_0^\infty K(x; u; h_n^*) dF(u) \right| \\ &= \left| \int_0^\infty K(x; u; h_n^*) d[F_n(u) - F(u)] \right| \\ &= \left| \int_0^\infty [F_n(u) - F(u)] dK_u(x; u; h_n^*) \right| \\ &\leq \sup_{0 < t \leq T} |F_n(t) - F(t)| \left| \int_0^\infty dK_u(x; u; h_n^*) \right|. \end{aligned} \quad (13)$$

از طرفی بنا بر قضیه ۲ هوروات (۱۹۸۵)، F_n برآوردگر به طور یکنواخت سازگار قوی برای $F(t)$ است. بنابراین هرگاه $n \rightarrow \infty$ داریم

$$\sup_{0 < t \leq T} |F_n(t) - F(t)| \xrightarrow{a.s.} 0.$$

همچنین برای n های بزرگ

$$\left| \int_0^{\infty} dK_u(x; u, h_n^*) \right| = \frac{1}{\sqrt{2\pi h_n^* x}} e^{-\frac{1}{2h_n^* x}}.$$

بنابراین، وقتی $n \rightarrow \infty$ ، آنگاه $J_1 \rightarrow 0$. از طرف دیگر، مشابه برهان قضیه ۲ کولاسکرا و پاچت (۲۰۰۶)، داریم $J_2 \rightarrow 0$. پس اثبات کامل می شود. □

در ادامه پهنای باند برآوردگر جونز، بر اساس تابع هسته گاوسی و چگالی پیشین گاما وارونه با پارامترهای α و β به صورت

$$\tau(h) = \frac{2}{\Gamma(\alpha)\beta^\alpha} \frac{1}{h^{2\alpha+1}} \exp\left(-\frac{1}{\beta h^2}\right), \quad \alpha > 0, \beta > 0, h > 0, \quad (14)$$

محاسبه می شود.

قضیه ۳: برای تابع هسته گاوسی، توزیع پیشین (۱۴) و تابع زیان توان دوم خطا، برآوردگر بیزی h در نقطه x عبارت است از

$$h_n^* = \frac{\Gamma(\alpha) \sum_{i=1}^n Y_i^{-1} \{1/(\beta(Y_i - x)^2 + 2)\}^\alpha}{\sqrt{2\beta} \Gamma(\alpha + \frac{1}{2}) \sum_{i=1}^n Y_i^{-1} \{1/(\beta(Y_i - x)^2 + 2)\}^{\alpha + \frac{1}{2}}}. \quad (15)$$

برهان: با در نظر گرفتن رابطه (۷) و جایگذاری تابع هسته گاوسی در رابطه (۶) داریم

$$\begin{aligned} \hat{\pi}(h|Y_1, \dots, Y_n, x) &= \frac{\hat{f}_h(x)\pi(h)}{\int \hat{f}_h(x)\pi(h)dh} \\ &= \frac{\sum_{i=1}^n Y_i^{-1} \frac{1}{h^{2\alpha+2}} \exp\{-\frac{1}{h^2} \frac{(Y_i-x)^2}{2} + \frac{1}{\beta}\}}{\int_0^{\infty} \sum_{i=1}^n Y_i^{-1} \frac{1}{h^{2\alpha+2}} \exp\{-\frac{1}{h^2} \frac{(Y_i-x)^2}{2} + \frac{1}{\beta}\} dh} \end{aligned} \quad (16)$$

فرض کنید $\mu_i = \frac{1}{2}(Y_i - x)^2 + \frac{1}{\beta}$ ، طرف راست رابطه (۱۶) به صورت

$$= \frac{\sum_{i=1}^n Y_i^{-1} \left(\frac{1}{h^{2\alpha+2}}\right) \exp\left\{-\frac{1}{h^2} \frac{(Y_i-x)^2}{2} + \frac{1}{\beta}\right\}}{\sum_{i=1}^n Y_i^{-1} \int_0^{\infty} \exp\left(-\frac{\mu_i}{h^2}\right) \left(\frac{1}{h^{2\alpha+2}}\right) dh} \quad (17)$$

۵۰..... برآورد بیزی و اعتبار سنجی متقابل پهنای باند برآوردگر هسته‌ای تابع چگالی

حاصل می‌شود. برای بررسی منجر کسر (۱۷) فرض کنید

$$u = \frac{1}{h^2}$$

داریم

$$\begin{aligned} \sum_{i=1}^n Y_i^{-1} \int_0^{\infty} \exp\left(-\frac{\mu_i}{h^2}\right) \frac{1}{h^{2\alpha+2}} dh &= -\frac{1}{2} \sum_{i=1}^n Y_i^{-1} \int_0^{\infty} \exp(-\mu_i u) u^{\alpha-\frac{1}{2}} du \\ &= \frac{1}{2} \sum_{i=1}^n Y_i^{-1} \Gamma\left(\alpha + \frac{1}{2}\right) \mu_i^{-(\alpha+\frac{1}{2})} \\ &= \sum_{i=1}^n Y_i^{-1} \frac{\Gamma\left(\alpha + \frac{1}{2}\right)}{2} \left\{ \frac{(Y_i - x)^2}{2} + \frac{1}{\beta} \right\}^{-\alpha-\frac{1}{2}} \end{aligned}$$

بنابراین

$$\hat{\pi}(h|Y_1, \dots, Y_n, x) = \frac{\sum_{i=1}^n Y_i^{-1} \left(\frac{1}{h^{2\alpha+2}}\right) \exp\left\{-\frac{1}{h^2} \frac{(Y_i - x)^2}{2} + \frac{1}{\beta}\right\}}{\sum_{i=1}^n Y_i^{-1} \left(\Gamma\left(\alpha + \frac{1}{2}\right)/2\right) \left\{(Y_i - x)^2/2 + \frac{1}{\beta}\right\}^{-\alpha-\frac{1}{2}}}$$

پهنای باند موضعی به صورت

$$h_n^*(x) = \frac{\int_0^{\infty} \sum_{i=1}^n Y_i^{-1} \exp(-\mu_i/h^2) \left(\frac{1}{h^{2\alpha+2}}\right) dh}{\sum_{i=1}^n Y_i^{-1} \left(\Gamma\left(\alpha + \frac{1}{2}\right)/2\right) \left\{(Y_i - x)^2/2 + \frac{1}{\beta}\right\}^{-\alpha-\frac{1}{2}}}$$

محاسبه می‌شود، که صورت کسر آن عبارت است از

$$\begin{aligned} \int_0^{\infty} \sum_{i=1}^n Y_i^{-1} \exp\left(\frac{-\mu_i}{h^2}\right) \frac{1}{h^{2\alpha+2}} h dh &= \frac{1}{2} \sum_{i=1}^n Y_i^{-1} \int_0^{\infty} \exp(-\mu_i u) u^{\alpha-1} du \\ &= \frac{1}{2} \sum_{i=1}^n Y_i^{-1} \Gamma(\alpha) \mu_i^{-\alpha} \\ &= \frac{\Gamma(\alpha)}{2} \sum_{i=1}^n Y_i^{-1} \left\{ \frac{(Y_i - x)^2}{2} + \frac{1}{\beta} \right\}^{-\alpha} \end{aligned}$$

بنابراین

$$h_n^*(x) = \frac{\Gamma(\alpha) \sum_{i=1}^n Y_i^{-1} \left\{(Y_i - x)^2/2 + \frac{1}{\beta}\right\}^{-\alpha}}{\Gamma\left(\alpha + \frac{1}{2}\right) \sum_{i=1}^n Y_i^{-1} \left\{(Y_i - x)^2/2 + \frac{1}{\beta}\right\}^{-\alpha-\frac{1}{2}}}$$

$$= \frac{\Gamma(\alpha) \sum_{i=1}^n Y_i^{-1} \{1/(\beta(Y_i - x)^2 + 2)\}^\alpha}{\sqrt{2\beta} \Gamma(\alpha + \frac{1}{2}) \sum_{i=1}^n Y_i^{-1} \{1/(\beta(Y_i - x)^2 + 2)\}^{\alpha + \frac{1}{2}}}$$

تذکر ۱: برای بررسی سازگاری قوی برآوردگر $\hat{f}_{h_n^*}$ با هسته گاوسی می توان مشابه قضیه ۲ عمل کرد.

۳ شبیه سازی

در این بخش، برای مقایسه برآورد پهنای باند به دو روش بیزی و اعتبارسنجی متقابل کمترین توان های دوم مطالعه ای شبیه سازی با دو هسته گاوسی و گاوسی وارون انجام می شود. در این شبیه سازی ها، توزیع f عضو خانواده های گاما به صورت

$$f(t) = \frac{t^{\gamma-1} e^{-\frac{t}{\theta}}}{\theta^\gamma \Gamma(\gamma)}, \gamma > 0, \theta > 0, t \geq 0,$$

و وایبول به صورت

$$f(t) = \frac{\gamma t^{\gamma-1} e^{-(\frac{t}{\theta})^\gamma}}{\theta^\gamma}, \gamma > 0, \theta > 0, t \geq 0,$$

در نظر گرفته می شود. طبق رابطه (۴)، اگر توزیع نااریب $Gamma(\gamma, \theta)$ باشد، آنگاه توزیع جامعه در طول-اریب $Gamma(\gamma + 1, \theta)$ خواهد بود. به طریق مشابه اگر توزیع نااریب $Weibull(\gamma, \theta)$ باشد، آنگاه توزیع جامعه در طول-اریب گاما تعمیم یافته به صورت

$$f(t) = \frac{\gamma \theta^{-\gamma(1+\frac{1}{\gamma})} t^{\gamma(1+\frac{1}{\gamma})-1} e^{-(\frac{t}{\theta})^\gamma}}{\Gamma(1+\frac{1}{\gamma})}, \gamma > 0, \theta > 0, t \geq 0.$$

است. در محاسبه برآورد بیزی h از توزیع های پیشین (۹) و (۱۴) به ترتیب برای هسته های گاوسی وارون و گاوسی استفاده می شود. در این شبیه سازی ها اثر پارامترهای توزیع پیشین روی برآوردگر $\hat{f}_{h^*}(t)$ بررسی شده است.

۵۲ برآورد بیزی و اعتبارسنجی متقابل پهنای باند برآوردگر هسته‌ای تابع چگالی

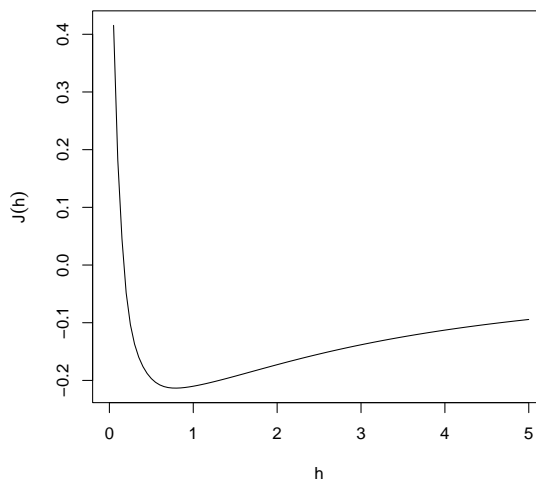
گیامون و همکاران (۱۹۹۸) نشان دادند در روش اعتبارسنجی متقابل کمترین توان‌های دوم، وقتی هسته متقارن است، h بهینه را می‌توان از کمینه کردن عبارت

$$J(h) \cong \left(\sum_{i=1}^n y_i^{-1} \right)^{-2} h^{-1} \sum_{j=1}^n \sum_{i=1}^n (y_i y_j)^{-1} K * K \left(\frac{y_i - y_j}{h} \right) - 2 \mu n^{-1} \sum_{i=1}^n \left(\sum_{j \neq i} y_j^{-1} \right)^{-1} y_i^{-1} h^{-1} \sum_{j \neq i} y_j^{-1} K \left(\frac{y_j - y_i}{h} \right),$$

به دست آورد، که در آن $K * K(\cdot)$ پیچش دو تابع K است و به صورت زیر تعریف می‌شود

$$K * K(u) = \int_{-\infty}^{\infty} K(u-v)K(v)dv.$$

نمودار $J(h)$ در شکل ۱، رسم شده است. به طریق مشابه، اگر هسته نامتقارن باشد،



شکل ۱: نمودار $J(h)$ در جامعه ناریب توزیع گاما با پارامترهای $\alpha = \beta = 3$ و $n = 20$ ، $\theta = 1.7 = 2$

می‌توان با جایگذاری رابطه (۶) در عبارت زیر و استفاده از روش اعتبارسنجی

متقابل کمترین توان‌های دوم، h بهینه را به دست آورد

$$\hat{J}(h) \cong \int \hat{f}_h^2(t) dt - 2 \int \hat{f}_h(t) f(t) dt.$$

برای مقایسه عملکرد دو روش بیزی و اعتبارسنجی متقابل در برآورد h ، نسبت MSE های برآورد شده برآوردگر جونز با دو پهنای باند به ازای مقادیر مختلف t به صورت

$$r(t) = \frac{EMSE(\hat{f}_{h_c}(t))}{EMSE(\hat{f}_{h_n^*}(t))},$$

محاسبه و نمودار آن در شکل‌های ۲ و ۴ رسم شده است، که در آن

$$EMSE(\phi_n(t)) = \frac{1}{N} \sum_{i=1}^N (\phi_n(t) - \phi(t))^2.$$

$\phi(t)$ و $\phi_n(t)$ به ترتیب نشانگر تابع چگالی دلخواه و برآوردگر آن هستند و N تعداد دفعات شبیه سازی‌ها برابر ۱۰۰۰ در نظر گرفته شده است. h_c و h_n^* به ترتیب نشانگر مقدار h به دست آمده از روشهای اعتبارسنجی متقابل و بیزی هستند. بدیهی است که هر چقدر $EMSE(\hat{f}_{h_n^*}(t))$ کوچکتر از $EMSE(\hat{f}_{h_c}(t))$ باشد، روش بیزی از روش اعتبارسنجی متقابل بهتر عمل می‌کند. به عبارت دیگر برآوردگر به دست آمده با پهنای نوار بیزی بهتر از برآوردگر با پهنای نوار اعتبارسنجی متقابل است.

در شکل‌های ۲ و ۳ توزیع جامعه نارایب، گاما با پارامترهای $\theta = 1$ و $\gamma = 2$ بوده و هسته گاوسی است. همچنین پارامترهای توزیع پیشین h ، $\alpha = \beta = 3$ بوده و حجم نمونه $n = 20$ است.

با توجه به شکل ۲، در نزدیکی مبدأ مختصات، ملاحظه می‌شود که روش بیزی در برآورد پهنای باند بهتر از روش اعتبارسنجی متقابل عمل می‌کند. اما به تدریج با زیاد شدن t ، ابتدا روش اعتبارسنجی متقابل بهتر عمل کرده ولی با افزایش t هر دو روش تقریباً شبیه هم عمل می‌کنند.

شکل ۳، نمودار تابع چگالی جامعه نارایب گاما را به همراه برآوردگر جونز با پهنای باندهای برآورد شده به دو روش بیزی و اعتبارسنجی متقابل را نشان می‌دهد. همانگونه که در این شکل مشاهده می‌شود برآورد تابع چگالی به روش بیزی تا

۵۴ برآورد بیزی و اعتبارسنجی متقابل پهنای باند برآوردگر هسته‌ای تابع چگالی

نزدیکی نقطه ۲ نزدیک‌تر به تابع چگالی اصلی می‌باشد. بین نقاط ۲ تا ۵ برآورد تابع چگالی به روش اعتبارسنجی متقابل بهتر عمل می‌کند. از نقطه ۵ به بعد دو روش تقریباً شبیه هم عمل می‌کنند.

در شکل‌های ۴ و ۵ توزیع جامعه نارایب، وایبول با پارامترهای $\theta = 1$ و $\gamma = 2$ بوده و هسته گاوسی می‌باشد. همچنین پارامترهای توزیع پیشین h ، $\alpha = \beta = 3$ بوده و حجم نمونه $n = 20$ است. با توجه به شکل ۴، در نزدیکی مبدأ مختصات، ملاحظه می‌شود که روش بیزی در برآورد پهنای باند بهتر از روش اعتبارسنجی متقابل عمل می‌کند. اما به تدریج با زیاد شدن t ، روش اعتبارسنجی متقابل بهتر عمل می‌کند.

شکل ۵، نمودار تابع چگالی توزیع نارایب وایبول را به همراه نمودارهای برآوردگر جونز با پهنای باندهای برآورد شده به دو روش بیزی و اعتبارسنجی متقابل نشان می‌دهد. با توجه به این نمودار دو روش بیزی و اعتبارسنجی متقابل تقریباً شبیه هم عمل می‌کنند.

در ادامه ملاحظه می‌شود چگونه انتخاب پارامترهای توزیع جامعه نارایب (β) و (α) و توزیع پیشین h یعنی (θ و γ) و همچنین حجم نمونه گیری دارای نقش مهمی در ساختار برآوردگرهای جونز هستند.

در جدول ۱، مقدار جمع بسته میانگین توان دوم خطای برآورد شده برآوردگرها (EIMSE) با در نظر گرفتن هسته گاوسی و حجم‌های نمونه مختلف محاسبه شده است. در توزیع پیشین (۱۴)، برای مقادیر متفاوت β ، α برابر ۳ در نظر گرفته شده است. این نسبت به صورت

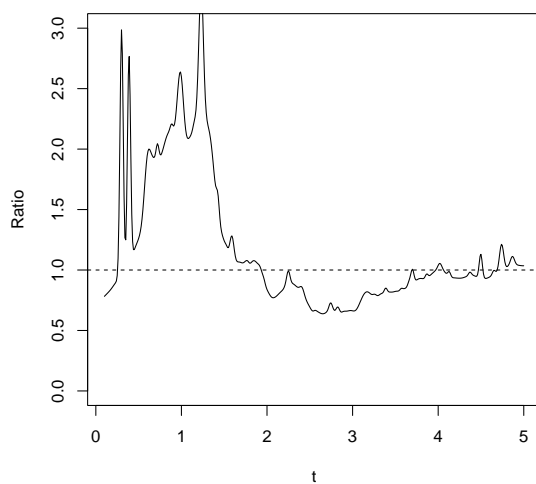
$$r^* = \frac{EIMSE(\hat{f}_{hc})}{EIMSE(\hat{f}_{h^*})}$$

است، که در آن

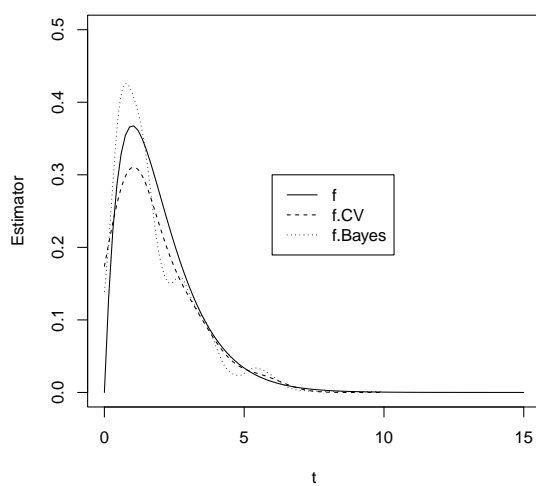
$$EIMSE(\phi_n) = \frac{1}{mN} \sum_{j=1}^N \sum_{i=1}^m (\phi_n(z_i) - \phi(z_i))^2.$$

برای محاسبه مقدار فوق، z_i ها به صورت

$$z_i = Y_{(1)} + \frac{i}{m}(Y_{(n)} - Y_{(1)}), \quad i = 1, \dots, m$$

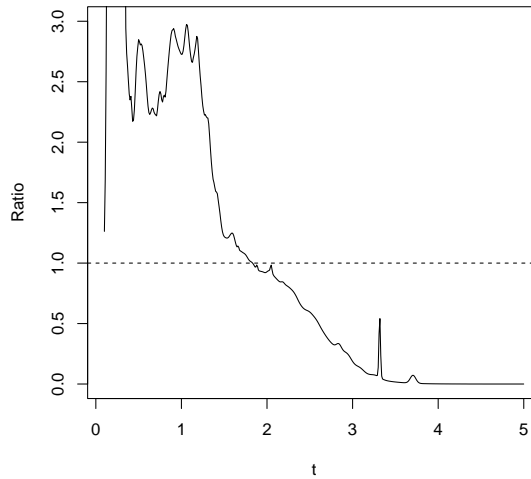


شکل ۲: نمودار $r(t)$ در خانواده گاما با پارامترهای $\alpha = \beta = 3, \theta = 1, \gamma = 2$ و $n = 20$

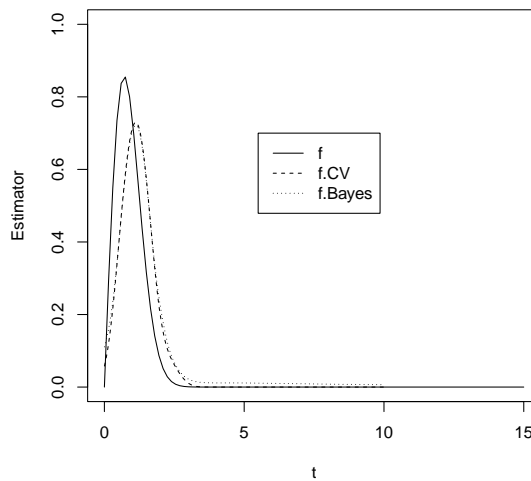


شکل ۳: نمودار تابع چگالی گاما و برآوردگر جونز با پهنای باندهای برآورد شده به دو روش بیزی و اعتبارسنجی متقابل

۵۶..... برآورد بیزی و اعتبارسنجی متقابل پهنای باند برآوردگر هسته‌ای تابع چگالی



شکل ۴: نمودار $r(t)$ در خانواده نارایب و ایبول با پارامترهای $\alpha = \beta = 3$ ، $n = 20$ و $\gamma = 1$ ، $\theta = 2$



شکل ۵: نمودار تابع چگالی و ایبول و برآوردگر جونز با پهنای باندهای برآورد شده به دو روش بیزی و اعتبارسنجی متقابل

مسعود عجمی بختیاروند، وحید فکور، سارا جمهوری ۵۷

انتخاب می‌شوند، که در آن $Y_{(j)}$ ، j -امین آماره مرتب در نمونه تصادفی Y_1, \dots, Y_n است.

جدول ۱: نسبت r^* برآوردگر جونز با پهنای باندهای اعتبارسنجی متقابل و بیزی در توزیع گاما و وایبول با $\theta = 1$ و γ مختلف

$Weibull(2, 1)$	$Gamma(2, 1)$	$Exp(1)$	$Weibull(0.5, 1)$	$Gamma(0.5, 1)$	β	n
۲/۹۴	۱/۹۹	۲/۷۵	۰/۸۰	۱/۹۸	۳	۵
۱/۵۹	۱/۳۲	۱/۹۰	۰/۵۵	۱/۵۲	۵	
۱/۰۵	۱/۲۱	۱/۷۵	۰/۲۹	۱/۰۷	۱۰	
۱/۶۶	۱/۴۳	۳/۳۵	۰/۹۸	۲/۱۱	۳	۲۰
۱/۷۴	۱/۰۱	۲/۶۷	۰/۸۳	۲/۱۶	۵	
۱/۲۶	۰/۶۶	۱/۷۳	۰/۶۸	۱/۹۱	۱۰	
۱/۸۵	۱/۵۶	۳/۲۱	۱/۰۵	۲/۱۶	۳	۴۰
۲/۳۸	۱/۱۵	۳/۰۹	۰/۹۱	۲/۴۹	۵	
۲/۰۹	۰/۷۷	۲/۸۵	۰/۷۶	۲/۶۱	۱۰	

همان‌طور که در جدول ۱ ملاحظه می‌شود که در تمامی حجم نمونه‌های انتخابی با افزایش پارامتر β مقدار نسبت r^* کاهش می‌یابد. لذا روش بیزی برای انتخاب پهنای باند در برآوردگر جونز با داده‌های در طول-اریب برای β های بزرگ بدتر از روش اعتبارسنجی متقابل عمل می‌کند. همچنین در خانواده وایبول در حجم نمونه پایین روش اعتبارسنجی متقابل بهتر از روش بیزی است.

۴ بحث و نتیجه‌گیری

در این مقاله روش بیزی برای انتخاب پهنای باند در برآوردگر جونز با داده‌های در طول-اریب استفاده شد و با روش اعتبارسنجی متقابل مقایسه گردید. قابل ذکر است که روش بیزی معایبی نیز دارد که از آن جمله می‌توان به مشکل انتخاب یک توزیع پیشین مناسب اشاره کرد. در اینجا انتخاب توزیع پیشین وابسته به نوع تابع هسته به کار رفته در محاسبه برآوردگر است.

مواردی وجود دارند که نیازمند تحقیق و بررسی بیشتر هستند و آن استفاده از توزیع‌های پیشین مناسب دیگر و هسته‌های متقارن و نامتقارن دیگر است. علاوه بر این همگرایی یکنواخت قوی برآوردگر تابع چگالی جونز و همگرایی میانگین مرتبه

۵۸ برآورد بیزی و اعتبار سنجی متقابل پهنای باند برآوردگر هسته‌ای تابع چگالی

دوم آن نیز برای خانواده هسته‌های متقارن و نامتقارن دیگر قابل بررسی است. علاقمندان برای دریافت برنامه‌های این مقاله که در محیط نرم افزار R تهیه شده‌اند می‌توانند با آدرس الکترونیک نویسنده مسئول مقاله تماس حاصل نمایند.

تقدیر و تشکر

نویسندگان از پیشنهادات داوران و هیئت تحریریه محترم مجله که باعث اصلاحات سازنده در این مقاله شده‌اند، کمال تشکر را دارند.

مراجع

- Bhattacharyya, B. B., Franklin, L. A. and Richardson, G. D. (1988), A Comparison of Nonparametric Unweighted and Length-biased Density Estimation of Fibres, *Communications in Statistics-Theory and Methods*, **17**, 3629-3644.
- Blumenthal, S. (1967), Proportional Sampling in Life Length Studies, *Technometrics*, **9**, 205-218.
- Chaubey, Y. P., Sen, P. K. and Li, J. (2010), Smooth Density Estimation for Length-biased Data, *Journal of the Indian Society of Agricultural Statistics*, **64**, 145-155.
- Cox, D. R. (1969), Some Sampling Problems in Technology, *New Developments in Survey Sampling*, Edited by Johnson and Smith, Wiley.
- Cristobal, J. A. and Alcalá, J. T. (2001), An Overview of Nonparametric Contributions to the Problem of Functional Estimation from Biased Data, *Sociedad de Estadística e Investigación Operativa*, **10**, 309-332.

۵۹..... مسعود عجمی بختیاروند، وحید فکور، سارا جمهوری

Efromovich, S. (2004), Density Estimation for Biased Data, *The Annals of Statistics*, **32**, 1137-1161.

Gangopadhyay, A. K. and Cheung, K. N. (2002), Bayesian Approach to the Choice of Smoothing Parameter in Kernel Density Estimation, *Nonparametric Statistics*, **14**, 655-664.

Guillamon, A., Navarro, j. and Ruiz, J. M. (1998), Kernel Density Estimation Using Weighted Data, *Communications in Statistics-Theory and Methods*, **27**, 2123-2135.

Horváth, L. (1985), Estimation From a Length-Biased Distribution, *Statistics and Decisions*, **3**, 91-113.

Jones, M. C. (1991), Kernel Density Estimation for Length Biased Data, *Biometrika* **78**, 511-519.

Kulasekera, K. B. and Padgett, W. J. (2006), Bayes Bandwidth Selection in Kernel Density Estimation with Censored Data, *Nonparametric Statistics*, **18**, 129-143.

Mcfadden, J. A. (1962), On the Lengths of Intervals in a Stationary Point Process, *Journal of the Royal Statistical Society*, **B 24**, 364-382.

Parzen, E. (1962), On Estimation of a Probability Density Function and Mode, *The Annals of Mathematical Statistics*, **33**, 1065-1076.

Rosenblatt, M. (1956), Remarks on Some Nonparametric Estimation of a Density Function, *The Annals of Mathematical Statistics*, **27**, 832-837.

Silverman, B. W. (1985), *Density Estimation for Statistics and Data Analysis*, Chapman and Hall, London.

۶۰..... برآورد بیزی و اعتبار سنجی متقابل پهنای باند برآوردگر هسته‌ای تابع چگالی

Wand, M. P and Jones, M. C.(1985), *Kernel Smoothing*, Chapman and Hall, London.

Wicksell, S. D. (1925), The Corpuscle Problem. A Mathematical Study of a Biometrica Problem, *Biometrika*, **17**, 84-99.