

## به دست آوردن برآوردگرهای گشتاوری با استفاده از تابع توزیع تجربی

بهزاد منصوری<sup>۱</sup>، رحیم چینی پرداز<sup>۲</sup> سامی عطیه سید الفرطوسی<sup>۳</sup> و حبیب اله ممبینی<sup>۴</sup>

تاریخ دریافت: ۱۴۰۱/۰۵/۰۸

تاریخ پذیرش: ۱۴۰۱/۱۲/۰۲

### چکیده:

از تابع توزیع تجربی به عنوان برآورد تابع توزیع احتمال تجمعی یک متغیر تصادفی استفاده می‌شود. تابع توزیع تجربی نقشی اساسی در بسیاری از استنباط‌های آماری دارد که در برخی موارد کمتر شناخته شده هستند. در این مقاله تابع احتمال تجربی به عنوان مشتق تابع توزیع تجربی معرفی شده و با استفاده از آن نشان داده می‌شود که برآوردگرهای گشتاوری مانند میانگین نمونه، میانه نمونه، واریانس نمونه‌ای و ضریب همبستگی نمونه‌ای حاصل جایگزین کردن تابع چگالی متغیر تصادفی با تابع احتمال تجربی آن در تعریف‌های نظری هستند. علاوه بر این، برای برآورد پارامترهای جامعه از برآوردگر هسته تابع چگالی احتمال استفاده شده و یک روش جدید برای برآورد پهنای باند در برآوردگر هسته تابع چگالی احتمال، معرفی شده است.

واژه‌های کلیدی: تابع توزیع تجربی، برآوردگر هسته، برآورد گشتاوری، پهنای باند

### ۱ مقدمه

حمله قلبی ( $p$ ) به صورت زیر محاسبه می‌شود:

$\hat{p} = \frac{1}{n} \{ \text{تعدادی از } n \text{ نفر که } CRP \text{ آن‌ها کمتر یا مساوی } 3 \text{ است} \}$ .

به سادگی می‌توان نشان داد که  $\hat{p} \xrightarrow{P} p$  هرگاه  $n \rightarrow \infty$ . در اینجا  $\hat{p}$  نماد همگرایی در احتمال است.

در مثال بالا فرض کنید که متغیر تصادفی  $X$  نشان‌دهنده میزان  $CRP$  برای یک فرد معین در این جامعه باشد. در این صورت  $p = P(X \leq 3) = F_X(3)$  است که در آن  $F_X(x)$  تابع توزیع تجمعی متغیر تصادفی  $X$  است. مفهوم احتمال تجربی ارتباطی نزدیک با تابع توزیع تجربی دارد. برای نمونه تصادفی مستقل و هم‌توزیع  $X_1, \dots, X_n$  تابع توزیع تجربی به صورت

$$F_n(X) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$$

تعریف می‌شود که در آن  $I(x)$  تابع نشانگر است. در مثال بالا

داریم:

$$F_n(3) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq 3) = \hat{p}.$$

در علم آمار، برای احتمال رخداد یک پیشامد از یک نمونه تصادفی از فرآیند مولد آن پیشامد، کمک می‌گیرند و فراوانی نسبی حالت‌های مساعد در نمونه را برآوردی از احتمال واقعی، تلقی می‌کنند. در نظریه احتمال مقادیر اندازه‌گیری شده از نمونه تصادفی را اندازه‌گیری‌های تجربی و فراوانی نسبی آن‌ها را «احتمال تجربی» پیشامد می‌گویند. برای مثال در پزشکی  $CRP$ <sup>۵</sup> از جمله پروتئین‌هایی است که در هنگام وجود التهاب در بدن افزایش می‌یابد. از  $CRP$  می‌توان برای پی بردن به ریسک خطر ابتلا به بیماری‌های قلبی استفاده کرد. اگر مقدار  $CRP$  در خون بیشتر از ۳ باشد، پزشکان خطر ابتلا به بیماری قلبی عروقی را بسیار بالا ارزیابی می‌کنند. حال فرض کنید که علاقه‌مند هستیم نسبت افرادی که در یک جامعه مشخص کمتر در معرض بیماری‌های قلبی و عروقی هستند را برآورد کنیم. برای این منظور از  $CRP$  اندازه‌گیری شده در یک نمونه تصادفی از این جامعه استفاده می‌کنیم. با فرض اینکه حجم نمونه تصادفی  $n$  باشد، احتمال تجربی پایین بودن ریسک خطر

<sup>۱</sup> گروه آمار- دانشکده علوم ریاضی و کامپیوتر- دانشگاه شهید چمران اهواز (نویسنده مسئول: b.mansouri@scu.ac.ir)

<sup>۲</sup> گروه آمار- دانشکده علوم ریاضی و کامپیوتر- دانشگاه شهید چمران اهواز

<sup>۳</sup> گروه آمار- دانشکده علوم ریاضی و کامپیوتر- دانشگاه شهید چمران اهواز

<sup>۴</sup> گروه آمار- دانشکده علوم ریاضی و کامپیوتر- دانشگاه شهید چمران اهواز

## ۲ تابع توزیع تجربی

فرض کنید که  $X_1, \dots, X_n$  یک دنباله از متغیرهای تصادفی مستقل و هم‌توزیع با توزیع  $F(x)$  باشد. با توجه به هم‌توزیع و مستقل بودن  $X_i$  ها که دارای توزیع برنولی هستند،  $nF_n(x)$  دارای توزیع دوجمله‌ای با پارامترهای  $n$  و  $p = F(x)$  است؛ بنابراین

$$E(F_n(x)) = F(x), \text{Var}(F_n(x)) = \frac{F(x)(1-F(x))}{n}$$

در نتیجه  $MSE(F_n(x)) = n^{-1}F(x)(1-F(x))$ . بنابراین هرگاه  $n \rightarrow \infty$  آنگاه  $F_n(x) \xrightarrow{p} F(x)$ . مرتبه همگرایی تابع توزیع تجربی  $O_P(n^{-1/2})$  است یا به عبارتی با نرخ  $\sqrt{n}$  سازگار است زیرا برای  $\epsilon > 0$  دلخواه و ثابت متناهی  $M$  به صورت  $M = \frac{F(x)(1-F(x))}{\epsilon}$  از نامساوی چیبیشف داریم:

$$P(|\sqrt{n}(F_n(x) - F(x))| > M) \leq \frac{nE(F_n(x) - F(x))^2}{M^2} = \epsilon.$$

می‌توان با استفاده از قانون قوی اعداد بزرگ نشان داد که برای هر  $x$  ثابت  $F_n(x)$  با احتمال یک به  $F(x)$  همگرا است. به عبارت دیگر هرگاه  $n \rightarrow \infty$  آنگاه با احتمال یک داریم:

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x) \rightarrow E[I(X_i \leq x)] = F(x) \quad (1)$$

این نتیجه ثابت می‌کند که  $F_n(x)$  به صورت نقطه‌به‌نقطه  $F(x)$  همگرا است. [۳]، [۱] مستقلاً در سال ۱۹۳۳ ثابت کردند که  $F_n(x)$  به صورت یکنواخت نیز به  $F(x)$  همگرا است؛ یعنی هرگاه  $n \rightarrow \infty$  آنگاه با احتمال یک داریم:

$$\sup_x |F_n(x) - F(x)| \rightarrow 0.$$

دقت کنید که  $\sup_x |F_n(x) - F(x)|$  آماره کولموگروف-اسمیرنوف برای آزمون هم‌توزیعی است. علاوه بر این هرگاه  $n \rightarrow \infty$  از قضیه حد مرکزی داریم:

$$\sup_x \sqrt{n}|F_n(x) - F(x)| \sim N(0, F(x)(1-F(x))). \quad (2)$$

با استفاده از قانون قوی اعداد بزرگ  $F_n(x) - F(x)$  همگرا است درحالی‌که طبق قضیه حد مرکزی  $\sqrt{n}(F_n(x) - F(x))$  واگرا است (دامنه متغیر نرمال استاندارد کل اعداد حقیقی است). به بیان ساده‌تر  $\frac{1}{\sqrt{n}} \sum_{i=1}^n I(X_i \leq x)$  همگرا است اما  $\frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$  واگرای کجاست؟ حالت سؤال این است که مرز همگرایی و واگرایی کجاست؟ به عبارت دیگر تا چه توانی از  $n$  عبارت  $(F_n(x) - F(x))$  همگرا است؟

به همین دلیل تابع توزیع تجربی را برآوردگر فراوانی نسبی نیز می‌گویند. سادگی و نرخ همگرایی سریع تابع توزیع تجربی سبب رواج استفاده از آن شده است به گونه‌ای که در مباحث آمار ناپارامتری به صورت گسترده‌ای از تابع توزیع تجربی به عنوان برآورد تابع توزیع مجهول یک متغیر تصادفی استفاده می‌شود. در واقع بسیاری از استنباط‌های آماری بر پایه تابع توزیع تجربی هستند. برای مثال فلسفه روش بوت استرپ ناپارامتری در استفاده از نمونه‌گیری با جایگذاری برای باز نمونه‌گیری (دادن شانس برابر به مشاهدات) بر مبنای استفاده از تابع توزیع تجربی برای متغیر تصادفی تحت مطالعه است. از کاربردهای متنوع دیگر تابع توزیع تجربی، می‌توان به آزمون کولموگروف-اسمیرنوف در آزمون‌های هم‌توزیعی، برآورد تابع مفصل و برآورد منحنی عملکرد گیرنده (ROC)<sup>۶</sup> اشاره کرد. در استنباط آماری، یک روش رایج برآورد پارامترهای مجهول استفاده از روش گشتاوری است. علیرغم اینکه روش گشتاوری برآورد پارامتر، در کنار روش حداکثر درست‌نمایی از رایج‌ترین روش‌ها است، به ندرت به مبنای نظری آن توجه شده است. [۱۰] با تعریف برآورد جایگزین<sup>۷</sup> به رابطه بین تابع توزیع تجربی و برآورد گشتاوری اشاره کرده است.

در این مقاله هدف نشان دادن اهمیت تابع توزیع تجربی و نقش آن در به دست آوردن برآوردگرهای گشتاوری به صورت روشن و ساده است. برای این منظور با استفاده از تابع توزیع تجربی و تابع دلتای دیراک، تابع احتمال تجربی را تعریف کرده و نشان می‌دهیم که برآوردگرهای گشتاوری پارامترهای جامعه از جایگزین کردن تابع احتمال تجربی به جای تابع چگالی احتمال در تعریف‌های نظری به دست می‌آیند. علاوه بر این از برآوردگر هسته تابع چگالی احتمال برای برآورد پارامترهایی مانند میانگین و واریانس استفاده کرده و یک روش جدید برای برآورد پهنای باند پیشنهاد می‌کنیم. مقاله در پنج بخش تنظیم شده است. در بخش دوم برخی از خواص تابع توزیع تجربی به اختصار مرور شده است. در بخش سوم تابع احتمال تجربی معرفی شده است. در این بخش رابطه بین تابع توزیع تجربی و تابع احتمال تجربی با برآوردگرهای گشتاوری نشان داده شده است. در بخش چهارم مقاله با استفاده از برآوردگر هسته تابع چگالی احتمال، برآوردگرهایی برای میانگین و واریانس جامعه به دست آمده است. بخش پنجم مقاله به بحث و نتیجه‌گیری اختصاص یافته است.

<sup>6</sup> Receiver operating curve

<sup>7</sup> Plug-in estimator

<sup>8</sup> Pointwise convergence

سطح  $1 - \alpha$  برای  $F(x)$  با استفاده از نامساوی  $VC$  به دست می‌آید اما دامنه این فاصله اطمینان نسبت به فاصله اطمینانی که با استفاده از نامساوی  $DKW$  به دست آمد، بسیار بزرگ‌تر بوده و کارا نیست. در ادامه این دو فاصله اطمینان را برای یک مجموعه داده شبیه‌سازی شده مقایسه می‌کنیم.

داده‌های این قسمت یک مجموعه آمیخته از ۱۰۰۰ داده شبیه‌سازی شده از توزیع نرمال با میانگین ۱۰ و انحراف معیار ۱ و ۱۰۰۰ داده شبیه‌سازی شده از توزیع نرمال با میانگین ۱۵ و انحراف معیار ۱/۵ هستند. شکل (۲) چگالی این مجموعه داده و شکل (۳) تابع توزیع تجربی داده‌ها را به همراه فاصله اطمینان به دست آمده از نامساوی‌های  $DKW$  به رنگ قرمز و  $VC$  به رنگ آبی را نشان می‌دهد.

### ۳ تابع احتمال تجربی

در این بخش می‌خواهیم از تابع توزیع تجربی برای محاسبه احتمال در هر نقطه از دامنه متغیر تصادفی استفاده کنیم. برای این منظور از تابع دلتای دیراک استفاده کنیم. تابع دلتای دیراک را با  $\delta(x)$  نشان می‌دهیم. این تابع در خواص زیر صدق می‌کند:

$$\delta(x) = \begin{cases} 0 & x \neq 0, \\ \infty & x = 0 \end{cases}$$

و  $\int_{-\infty}^{\infty} \delta(x) dx = 1$  دقت کنید که

$$\int_{-\infty}^{\infty} \delta(x-a)g(x)dx = \int_{-\infty}^{\infty} \delta(x-a)[g(a) + (g(x) - g(a))]dx = g(a),$$

زیرا  $\delta(x-a)(g(x) - g(a)) = 0$ . حال تابع نشانگر را در نظر بگیرید:

$$I(x \geq 0) = \begin{cases} 1 & x \geq 0, \\ 0 & x < 0, \end{cases}$$

مشتق این تابع همه‌جا به جز در نقطه صفر برابر صفر است و چون در نقطه صفر ناپیوسته است مشتق‌پذیر نبوده و می‌توان مشتق آن را بینهایت دانست. به همین دلیل [۶] تابع دلتای دیراک را مشتق تابع نشانگر تشخیص دادند یعنی  $\delta(x) = \frac{dI(x \geq 0)}{dx}$ . بنابراین اگر تابع احتمال تجربی را به عنوان مشتق تابع توزیع تجربی تعریف کرده و با

قانون لگاریتم تکراری این مرز را مشخص می‌کند. بر طبق قانون لگاریتم تکراری<sup>۹</sup>

$$\lim_{n \rightarrow \infty} \frac{\sqrt{n} \sup_x |F_n(x) - F(x)|}{\sqrt{2 \ln \ln n}} \leq \frac{1}{2}.$$

برای توضیح بیشتر درباره اثر ضرایب  $\frac{1}{\sqrt{n}}$  و  $\frac{1}{\sqrt{2 \ln \ln n}}$  در شکل (۱) نمودار این ضرایب به عنوان توابعی از  $n$  رسم شده و سرعت همگرایی آن‌ها به صفر را نمایش داده شده است. دقت کنید که چگونه منحنی  $\frac{1}{\sqrt{2 \ln \ln n}}$  مقادیری بین منحنی  $\frac{1}{\sqrt{n}}$  و  $\frac{1}{n}$  قرار می‌گیرد. خینچین (۱۹۲۴)<sup>۱۰</sup> و کولموگروف (۱۹۲۹)<sup>۱۱</sup> قانون لگاریتم تکراری را معرفی و مطالعه کرده‌اند. برای مطالعه بیشتر در این زمینه [۴] فصل دوم را ببینید.

### ۱.۲ فاصله اطمینان برای تابع توزیع تجمعی

از رابطه (۲) نمی‌توان یک فاصله اطمینان کاربردی برای تابع توزیع مجهول  $F(x)$  به دست آورد زیرا حدود فاصله اطمینان به  $F(x)$  وابسته خواهند بود. [۲] نشان دادند که برای هر  $\epsilon > 0$

$$P(\sup_x |F_n(x) - F(x)| > \epsilon) \leq 2e^{-2n\epsilon^2}. \quad (3)$$

این نامساوی که به نامساوی دوورتسکی-کیفر-ولفویتز (DKW) معروف است امکان ساختن یک فاصله اطمینان کاربردی برای  $F(x)$  را فراهم می‌کند. برای این منظور کافی است که برای  $\epsilon_n = \sqrt{\frac{1}{2n} \ln\left(\frac{1}{\alpha}\right)}$

$$L(x) = \max\{F_n(x) - \epsilon_n, 0\}, U(x) = \min\{F_n(x) + \epsilon_n, 1\}, \quad (4)$$

آنگاه از رابطه (۳) داریم:

$$P(L(x) \leq F(x) \leq U(x)) \geq 1 - \alpha.$$

با استفاده از نظریه [۹] می‌توان نامساوی‌های کلی‌تری را به دست آورد ([۱۰]). در حالت یک‌بعدی، با استفاده از قضیه وینیک و چرونکیس (VC) می‌توان نشان داد که

$$P(\sup_x |F_n(x) - F(x)| > \epsilon) \leq \lambda(n+1)e^{-\frac{n\epsilon^2}{\lambda}}. \quad (5)$$

(برای جزئیات بیشتر [۱۰] را ببینید). با استفاده از (۵) و به ازای  $\epsilon_n = \sqrt{\frac{\lambda \ln\left(\frac{\lambda(n+1)}{\alpha}\right)}{n}}$  در حدود رابطه (۴) یک فاصله اطمینان در

<sup>9</sup>Law of the iterated logarithm

<sup>10</sup> Khinchin

<sup>11</sup> Kolmogorov

<sup>12</sup> Wasserman

$f_n(x)$  نشان دهیم، آنگاه داریم:

$$f_n(x) := \frac{dF_n(x)}{dx} = \frac{1}{n} \sum_{i=1}^n \delta(X_i - x).$$

توزیع تجربی دومتغیره به صورت

$$F_n(x, y) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x_i, Y_i \leq y_i)$$

تعریف می‌شود. تابع احتمال توأم تجربی دو متغیر تصادفی  $X$  و  $Y$  را به صورت زیر تعریف می‌کنیم،

$$f_n(x, y) := \frac{dF_n(x, y)}{dxdy} = \frac{1}{n} \sum_{i=1}^n \delta((X_i - x)(Y_i - y)).$$

امید ریاضی تجربی یک تابع مانند  $g(X, Y)$  عبارت است از

$$\begin{aligned} E_n(g(X, Y)) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_n(x, y) dx dy \\ &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) \delta((X_i - x)(Y_i - y)) dx dy \\ &= \frac{1}{n} \sum_{i=1}^n g(X_i, Y_i) \end{aligned}$$

بنابراین

$$E_n(XY) = \frac{1}{n} \sum_{i=1}^n X_i Y_i = \bar{XY},$$

در نتیجه

$$\sigma_n(xy) = E_n(XY) - E_n(X)E_n(Y) = \bar{XY} - \bar{X}\bar{Y},$$

برآوردگر تجربی کواریانس دو متغیر تصادفی است و برآوردگر

تجربی ضریب همبستگی نیز به صورت

$$r_n(x, y) = E_n(XY) = \frac{\sigma_n(x, y)}{\sigma_n(x)\sigma_n(y)} = \frac{\bar{XY} - \bar{X}\bar{Y}}{\sqrt{\bar{X}^2 - \bar{X}^2} \sqrt{\bar{Y}^2 - \bar{Y}^2}}$$

است؛ یعنی ضریب همبستگی نمونه‌های نیز حاصل جایگزین کردن توابع چگالی مجهول  $f(x)$ ،  $f(y)$  و  $f(x, y)$  با معادل‌های تجربی آن‌ها در تعریف نظری ضریب همبستگی است.

### ۲.۳ میانه تجربی

برای نمونه تصادفی مستقل و هم‌توزیع  $X_1, \dots, X_n$  فرض کنید که  $X_{(i)}, i = 1, \dots, n$  نشان‌دهنده  $i$  امین آماره ترتیبی در نمونه باشد. در این صورت اگر  $n$  فرد باشد آنگاه با تعریف

$$P_n(X \leq a) := F_n(a) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq a),$$

داریم:

$$P_n\left(X \leq X_{(\frac{n+1}{2})}\right) = P_n\left(X \geq X_{(\frac{n+1}{2})}\right),$$

یعنی  $X_{(\frac{n+1}{2})}$  میانه نمونه با استفاده از تابع توزیع تجربی داده‌ها است.

اگر  $n$  زوج باشد آنگاه با تعریف

$$Med = \frac{1}{2} \left[ X_{(\frac{n}{2})} + X_{(\frac{n}{2}+1)} \right],$$

در واقع تابع احتمال تجربی در همسایگی بسیار کوچک هر مشاهده وزن  $\frac{1}{n}$  را اعمال می‌کند. با وجود اینکه تابع احتمال تجربی به معنی دقیق یک تابع احتمال نیست، اما چنانکه در ادامه خواهیم دید، استفاده از آن برای نشان دادن فلسفه نظری برآوردگرهای گشتاوری و رابطه آن‌ها با تابع توزیع تجربی بسیار مفید است. با توجه به مستقل و هم‌توزیع بودن مشاهدات و تعریف تابع دلتای دیراک داریم:

$$E((f_n(x)) = E(\delta(Y - x)) = \int_{-\infty}^{\infty} \delta(Y - x) f(y) dy = f(x)$$

و

$$\begin{aligned} Var(f_n(x)) &= \frac{1}{n} Var(\delta(Y - x)) \\ &= \frac{1}{n} \{E(\delta(Y - x)^2) - (E(\delta(Y - x)))^2\} \\ &= \frac{f(x)(1 - f(x))}{n}. \end{aligned}$$

بنابراین هرگاه  $n \rightarrow \infty$  آنگاه  $f_n(x) \xrightarrow{P} f(x)$

اگر امید ریاضی تجربی یک تابع مانند  $g(x)$  را با  $E_n(g(X))$  نشان دهیم، با توجه به تعریف تابع دلتای دیراک داریم:

$$\begin{aligned} E_n(g(X)) &= \int_{-\infty}^{\infty} g(x) dF_n(x) \\ &= \int_{-\infty}^{\infty} g(x) f_n(x) dx \\ &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} g(x) \delta(X_i - x) dx \\ &= \frac{1}{n} \sum_{i=1}^n g(X_i). \end{aligned}$$

برای مثال

$$E_n(X) = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}, \quad E_n(X^2) = \frac{1}{n} \sum_{i=1}^n X_i^2 = \bar{X}^2$$

و

$$\sigma_n^2(X) = E_n(X^2) - (E_n(X))^2 = \bar{X}^2 - \bar{X}^2.$$

که عبارت فوق برآوردگر گشتاوری واریانس جامعه است.

### ۱.۳ تابع احتمال تجربی دومتغیره

برای نمونه تصادفی مستقل و هم‌توزیع  $(X_1, Y_1), \dots, (X_n, Y_n)$  با تابع توزیع احتمال توأم  $F(x, y)$  و تابع چگالی احتمال توأم  $f(x, y)$  تابع

داریم:

$$P_n(X \leq Med) = P_n(X \geq Med).$$

باید دقت کرد که همه برآوردگرهای نمونه‌ای حاصل جایگزین کردن تابع توزیع تجربی به جای تابع توزیع مجهول در روابط نظری نیستند. در بخش بعد نشان می‌دهیم که برخی از برآوردگرهای نارایب پارامترهای جامعه حاصل جایگزین کردن تابع توزیع مجهول با برآوردگر هستند.

#### ۴ برآوردگر چگالی هسته

در بخش قبل برآوردگرهای پارامترهای جامعه را با جایگزین کردن تابع احتمال تجربی به جای تابع احتمال مجهول به دست آوردیم. در این بخش می‌خواهیم برآوردگرها را با جایگزین کردن برآوردگر هسته تابع چگالی به جای تابع احتمال مجهول در روابط نظری به دست آورده و با نتایج بخش قبل مقایسه کنیم.

برای نمونه تصادفی مستقل و هم‌توزیع  $X_1, \dots, X_n$  با تابع چگالی احتمال  $f(x)$ ، برآوردگر تابع چگالی احتمال به فرم

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{x - X_i}{h}\right),$$

تعریف می‌شود که در آن  $h = h_n$  پهنای باند است و در شرایط  $\lim_{n \rightarrow \infty} nh = \infty$  و  $\lim_{n \rightarrow \infty} h = 0$  صدق می‌کند و تابع هسته  $k(x)$  نیز در شرایط

$$\int_{-\infty}^{\infty} k(x) dx = 1, \quad \int_{-\infty}^{\infty} xk(x) dx = 0, \\ \int_{-\infty}^{\infty} x^2 k(x) dx = \mu_2(k) \leq \infty, \quad (6)$$

صدق می‌کند. تحت این فرض‌ها و با فرض اینکه مشتق دوم تابع  $f(x)$  انتگرال‌پذیر است می‌توان نشان داد که  $\hat{f}(x)$  در احتمال به  $f(x)$  همگرا است ([۸] صفحات ۳۹-۴۰ را ببینید). همچنین می‌توان نشان داد که  $\hat{f}(x)$  با احتمال یک نیز به  $f(x)$  همگرا است ([۶] قضیه ۱-۴ و ۵-۱).

#### ۱.۴ امید ریاضی تابعی از یک متغیر تصادفی پیوسته با استفاده از برآوردگر چگالی هسته

فرض کنید  $g(x)$  که تابعی از متغیر تصادفی پیوسته  $X$  با تابع چگالی احتمال  $f(x)$  بوده و  $\hat{f}(x)$  برآوردگر هسته تابع چگالی با استفاده از

نمونه تصادفی  $X_1, \dots, X_n$  باشد. اگر  $E_{\hat{f}}$  نشان‌دهنده امید ریاضی با جایگزینی  $\hat{f}(x)$  به جای  $f(x)$  باشد و با فرض اینکه تابع هسته در روابط (۶) صدق کند داریم:

$$E_{\hat{f}}(g(X)) = \int_{-\infty}^{\infty} g(x) \hat{f}(x) dx \\ = \frac{1}{nh} \int_{-\infty}^{\infty} g(x) \sum_{i=1}^n k\left(\frac{x - X_i}{h}\right) dx \\ = \frac{1}{nh} \sum_{i=1}^n \int_{-\infty}^{\infty} g(x) k\left(\frac{x - X_i}{h}\right) dx \\ = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} g(t.h + X_i) k(t) dt \\ = \frac{1}{nh} \sum_{i=1}^n E_k(g(t.h + X_i)).$$

برای مثال

$$E_{\hat{f}}(X) = \frac{1}{n} \sum_{i=1}^n E_k(t.h + X_i) \\ = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} (t.h + X_i) k(t) dt = \bar{X},$$

زیرا طبق شرایط (۶) داریم:  $\int xk(x) dx = 0$ ،  $\int k(x) dx = 1$  همچنین

$$E_{\hat{f}}(X^2) = \frac{1}{n} \sum_{i=1}^n E_k(t.h + X_i)^2 \\ = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} (t^2 h^2 + 2t.hX_i + X_i^2) k(t) dt \\ = h^2 \int_{-\infty}^{\infty} t^2 k(t) dt + \bar{X}^2,$$

در این صورت:

$$Var_{\hat{f}}(\bar{X}) = E_{\hat{f}}(X^2) - (E_{\hat{f}}(X))^2 \\ = \bar{X}^2 - \bar{X}^2 + h^2 \mu_2(k),$$

که در آن  $\mu_2(k) = \int t^2 k(t) dt$  واریانس تابع هسته است. فرض کنید  $X_i$  ها نمونه‌ای تصادفی از توزیع  $f(x)$  با میانگین  $\mu$  و واریانس  $\sigma^2$  باشند. آنگاه با انتخاب  $\tilde{h} = \sqrt{\frac{\sigma^2}{n\mu_2(k)}}$  داریم:

$$Var_{\tilde{h}}(X) = E_{\tilde{h}}(X^2) - (E_{\tilde{h}}(X))^2 \\ = \bar{X}^2 - \bar{X}^2 + \frac{\sigma^2}{n},$$

و میدانیم که

$$E\left(\bar{X}^2 - \bar{X}^2 + \frac{\sigma^2}{n}\right) = \sigma^2,$$

یعنی با انتخاب درست پهنای باند،  $Var_{\tilde{h}}(X)$  برآوردگر واریانس با استفاده از تابع چگالی هسته، برآوردگر نارایب واریانس جامعه است درحالی‌که همان‌طور که در بخش قبل دیدیم  $\sigma_n^2$  (برآوردگر واریانس با

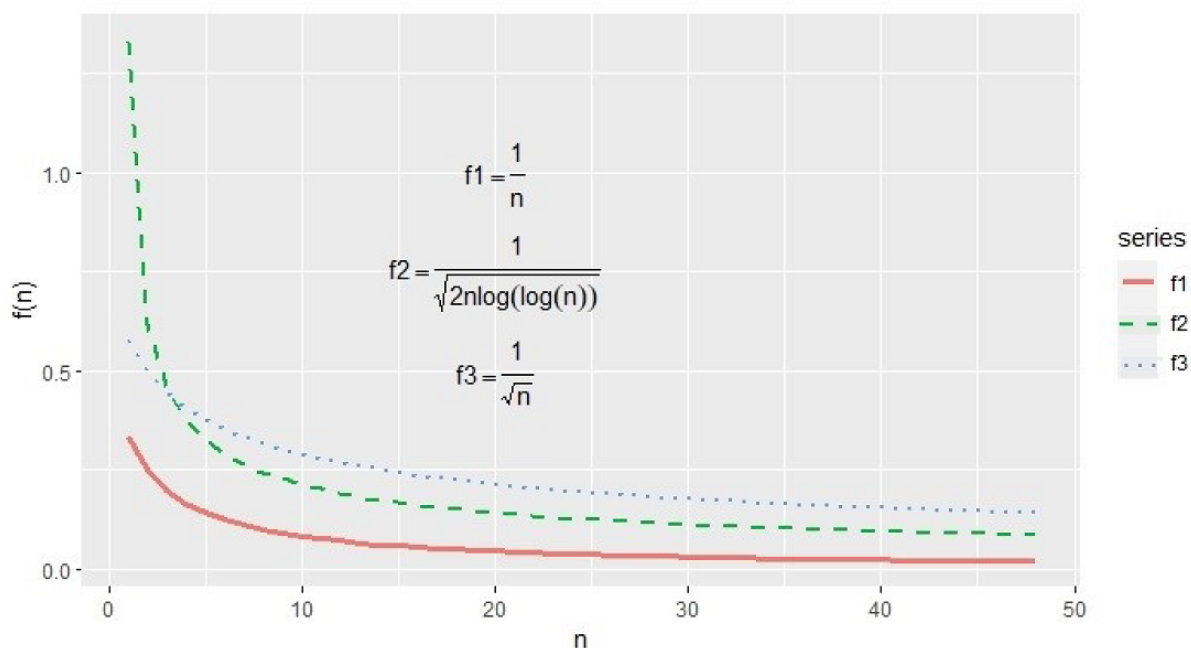
و نتیجه گرفت که اگر چگالی واقعی داده‌ها با توزیع نرمال متفاوت باشد مثلاً چند مدی، چوله یا برجسته باشد، کارایی این برآوردگر پایین بوده و سبب بیش همواری در چگالی برآورد شده می‌گردد. با توجه به اینکه  $\bar{h} < h_{ROT,f}$  انتظار این است که  $h_{ROT,f}$  نسبت به  $\bar{h}$  سبب بیش همواری در چگالی برآورد شده، شود. با توجه به سادگی  $\bar{h}$  مطالعه بیشتر آن و مقایسه عملکرد آن با دیگر برآوردگرهای پهنای باند می‌تواند موضوع جالبی برای تحقیق باشد.

استفاده از تابع توزیع تجربی) اریب است. دقت کنید که به ازای تابع هسته نرمال استاندارد  $\mu_x(k) = 1$  است و در نتیجه  $\bar{h} = \sigma n^{-\frac{1}{2}}$  به دست می‌آید.

[۸] با یک قاعده سرانگشتی<sup>۱۳</sup> برآوردگر زیر را برای پهنای باند به دست آورد:

$$h_{ROT,f} = 1/0.6\sigma n^{-\frac{1}{2}}$$

سیلورمن (صفحات ۴۸-۴۵) این برآوردگر را نقد و ارزیابی کرده



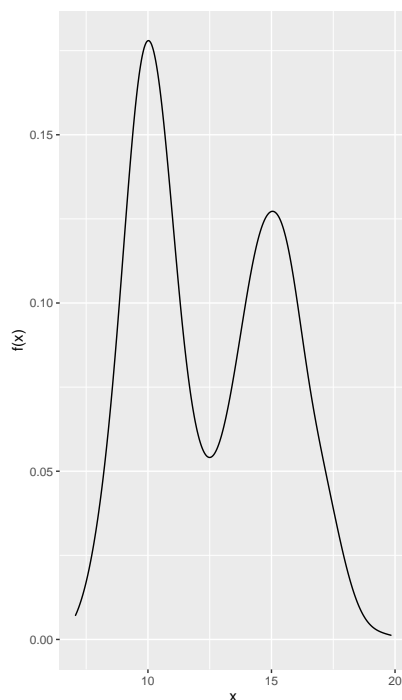
شکل ۱. مقایسه سرعت همگرایی دنباله  $\frac{1}{\sqrt{n}}$  و  $\frac{1}{\sqrt{2n \ln \ln n}}$  به ازای مقادیر  $n$  در فاصله ۳ تا ۵۰.

توزیع تجربی غیر هموار بودن آن است. به همین دلیل برآوردگرهای نوع هسته، رواج یافته‌اند. با این وجود برآوردگرهای نوع هسته نیز در برخی موارد تابع توزیع تجربی استفاده می‌کنند. برای مثال، در برآورد تابع توزیع احتمال به روش هسته، از تابع توزیع تجربی به‌عنوان یک محک استفاده می‌شود (به‌عنوان یک نمونه، [۷] را ببینید) با وجود قدمت تابع توزیع تجربی، تحقیق در مورد آن ادامه دارد ([۵]).

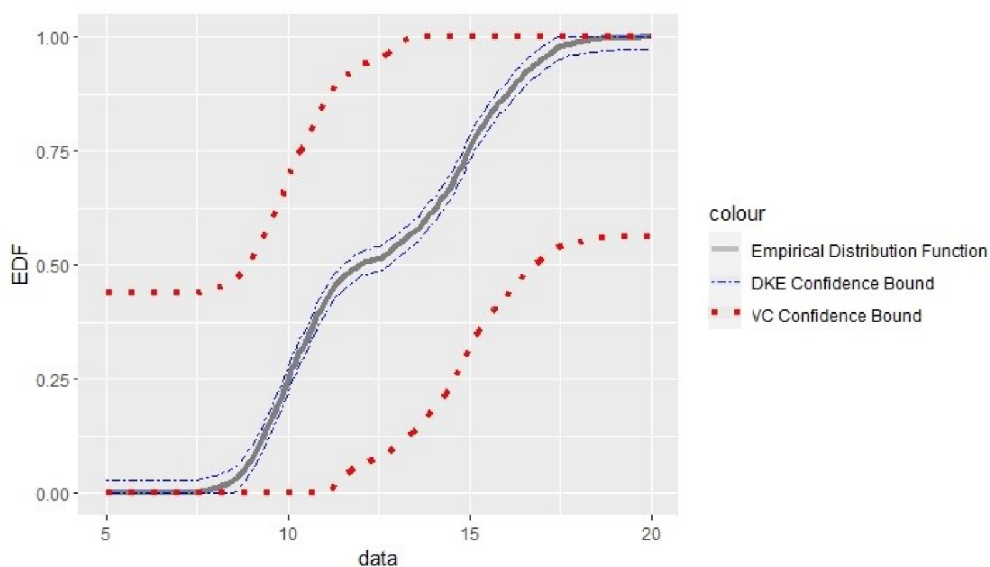
## ۵ بحث و نتیجه‌گیری

تابع توزیع تجربی، در عین سادگی نقشی پایه‌ای در آمار استنباطی به‌خصوص آمار ناپارامتری دارد. در این مقاله با استفاده از تابع توزیع تجربی و تابع دلتای دیراک، تابع احتمال تجربی را تعریف کرده و نشان دادیم که برآوردگرهای گشتاوری از جایگزین کردن تابع چگالی مجهول با تابع احتمال تجربی، در تعریف‌های نظری به دست می‌آیند. ضعف تابع

<sup>13</sup> Rule-of-thumb



شکل ۲. تابع چگالی نرمال آمیخته



شکل ۳. تابع توزیع تجربی داده‌ها (منحنی پیوسته) به همراه فاصله اطمینان به دست آمده از نامساوی‌های  $DKW$  منحنی خط-نقطه چین به رنگ آبی) و  $VC$  (منحنی نقطه چین به رنگ قرمز).

## مراجع

- [1] Cantelli, F. P. (1933). Sulla determinazione empirica delle leggi di probabilità. *Giornale dell'Istituto Italiano degli Attuari*, 4, 421-424.

- [2] Dvoretzky, A., Kiefer, J., and Wolfowitz, J. (1956). Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator. *The Annals of Mathematical Statistics*, **27(3)**, 642-669.
- [3] Glivenko, V. (1933). Sulla determinazione empirica delle leggi di probabilità. *Giornale dell'Istituto Italiano degli Attuari*, **4**, 92-99.
- [4] Kosorok, M. R. (2008). *Introduction to Empirical Processes and Semiparametric Inference*. New York: Springer.
- [5] Langrené, N., and Warin, X. (2021). Fast multivariate empirical cumulative distribution function with connection to kernel density estimation. *Computational Statistics and Data Analysis*, **162**, 107267.
- [6] Li, Q., and Racine, J. S. (2007). *Nonparametric Econometrics: Theory and Practice*. Princeton University Press.
- [7] Sarda, P. (1993). Smoothing parameter selection for smooth distribution functions. *Journal of Statistical Planning and Inference*, **35(1)**, 65-75.
- [8] Silverman, B. W. (1986). *Density Estimation for Statistics and Data*. Chapman and Hall: London
- [9] Vapnik, V. N., and Chervonenkis, A. Y. (1971). Theory of uniform convergence of frequency of appearance of attributes to their probabilities and problems of defining optimal solution by empiric data. *Avtomatika i Telemekhanika*, **2**, 42-53.
- [10] Wasserman, L. (2006). *All of Nonparametric Statistics*. Springer Science and Business Media. 13 Sarda.



## Obtaining moment estimators using the empirical distribution function

Behzad Mansouri<sup>1</sup>, Rahim Chinipardaz<sup>1</sup>, Sami Atiyah Sayyid Al-Farttosic<sup>1</sup> and Habiballah Mombeni<sup>1</sup>

### Abstract:

Empirical distribution function is used as an estimate of the cumulative probability distribution function of a random variable. The empirical distribution function has a fundamental role in many statistical inferences, which are little known in some cases. In this article, the empirical probability function is introduced as a derivative of the empirical distribution function, and it is shown that moment estimators such as sample mean, sample median, sample variance, and sample correlation coefficient result from replacing the random variable density function with the empirical probability function in the theoretical definitions. In addition, the kernel probability density function estimator is used to estimate the population parameters and a new method for bandwidth estimation in the kernel density estimation is introduced.

**Keywords:** Empirical distribution function, moment estimate, kernel estimator, bandwidth.

---

<sup>1</sup> Department of Statistics, Faculty of Mathematics and Computer Science, Shahid Chamran University of Ahvaz, Ahvaz, Iran