

کاربرد روش‌های خودگردان‌ساز در استنباط بیزی

مریم خان احمدی

دانشکده علوم ریاضی، دانشگاه شاهرود

چکیده

خودگردان‌ساز (بوت‌استرپ) از روش‌های باز نمونه‌گیری است که امروزه کاربرد گسترده‌ای در علم آمار دارد. این روش انواع مختلفی دارد. در این مقاله روش خودگردان‌ساز بیزی، که خود شامل روش‌های ناپارامتری و پارامتری است مورد بررسی قرار گرفته، و روش به‌دست آوردن آن‌ها در مثال‌های جداگانه ارائه شده است.

واژه‌های کلیدی: خودگردان‌ساز بیزی، استنباط ناپارامتری.

۱ تاریخچه

بعد از معرفی روش خودگردان توسط افرون، انواع مختلف خودگردان‌ساز پارامتری، خودگردان‌ساز هموار، خودگردان‌ساز دوگانه و خودگردان‌ساز بیزی معرفی شدند. خودگردان‌ساز بیزی برای اولین بار در سال ۱۹۸۱ توسط روبین^۴ [۱۵] معرفی شد و از آن برای برآورد توزیع پسین پارامتر مورد نظر استفاده کرد. لو^۵ [۱۰، ۱۱] این روش را به جوامع متناهی و به داده‌های سانسور شده تعمیم داد. خودگردان‌ساز بیزی و خودگردان‌ساز وزنی در مقاله‌های متعددی توسط: افرون [۵]، لو [۹]،

خودگردان‌ساز^۱ رده بزرگی از روش‌هاست، که در آن‌ها باز نمونه‌گیری^۲ از مجموعه داده‌های اولیه انجام شده و به همین دلیل روش‌های باز نمونه‌گیری نامیده می‌شوند. روش خودگردان‌ساز برای اولین بار توسط افرون^۳ [۵] در مقاله‌ای با عنوان «خودگردان‌ساز: نگاهی دیگر به روش جک‌نایف» معرفی و از آن برای برآورد اریبی، واریانس و توزیع نمونه‌ای آمارها استفاده شد. خودگردان‌ساز معرفی شده توسط افرون خودگردان‌ساز ناپارامتری بود.

^۱ Bootstrap

^۲ Resampling

^۳ Efron

^۴ Rubin

^۵ Lo

ونگ^۶ [۱۶]، میسون و نیوتن^۷ [۱۲]، پیریستگارد و

ولنر^۸ [۱۴]، گاسپارینی^۹ [۷] و جیمز^{۱۰} [۸] مورد

مطالعه قرار گرفت. بنکز^{۱۱} [۳] در مطالعات شبیه سازی از خودگردان ساز بیزی و نوع هموار آن استفاده کرد. خودگردان ساز بیزی ناپارامتری توسط روبین [۱۵]، اوون^{۱۲} [۱۳] مورد بررسی قرار گرفت.

توزیع پسین پارامتر شبیه سازی می شود. در استنباط بیزی فرض می شود θ متغیر تصادفی است که دارای یک توزیع پیشین است، و $\mathbf{X} = (X_1, \dots, X_n)$ یک نمونه تصادفی از F (توزیع شرطی X به شرط θ) است. استنباط بیزی بر اساس توزیع پسین $R(\mathbf{X}, F)$ به شرط X است. اگر F متعلق به یک خانواده پارامتری با پارامتر θ باشد، توزیع پسین $R(\mathbf{X}, F)$ در صورتی قابل محاسبه است که توزیع پیشین معلوم باشد.

۲ مقدمه

فرض کنید $\mathbf{X} = (X_1, \dots, X_n)$ یک نمونه تصادفی از توزیع نامعلوم F با پارامتر θ و $x = (x_1, \dots, x_n)$ یافته آن باشد. آماره $\hat{\theta} = \hat{\theta}(X)$ را به عنوان برآوردگر پارامتر θ در نظر می گیریم.

هدف برآورد توزیع نمونه ای آماره $\hat{\theta} = \hat{\theta}(X)$ می باشد. چون معمولاً از $\hat{\theta}$ برای ساختن فواصل اطمینان برای پارامتر نامعلوم θ (وابسته به F است) استفاده می کنیم به جای $\hat{\theta}$ یک ریشه آن یعنی $R_F = R(X_1, \dots, X_n, F)$ را در نظر می گیریم. ریشه R_F ممکن است به صورت $R_F = R(\mathbf{X}, F)$ باشد. پس نیاز داریم که توزیع نمونه ای $R_F = R(\mathbf{X}, F) = \hat{\theta} - \theta$ را به دست آوریم، که این کار با استفاده از روش خودگردان ساز ناپارامتری امکان پذیر است. با محاسبه آماره خودگردان ساز $\hat{\theta}^* = \hat{\theta}(X^*)$ می توان از توزیع $R(X^*, F_n) = \hat{\theta}^* - \hat{\theta}$ به عنوان تقریب توزیع نمونه ای

۳ خودگردان ساز بیزی ناپارامتری

در این بخش ابتدا الگوریتم محاسبه خودگردان ساز بیزی ناپارامتری را ارائه می کنیم. [۱] روبین [۱۵] یک روش خودگردان ساز بیزی ناپارامتری را در حالتی که توزیع پیشین نامعلوم است، ارائه کرد. هدف در این روش، پیدا کردن توزیع پسین $R(\mathbf{X}, F)$ است. ابتدا x_1, \dots, x_n به عنوان یک نمونه تصادفی مستقل و هم توزیع از متغیرهای تصادفی X_1, \dots, X_n با تابع توزیع F و تابع توزیع تجربی F_n را در نظر بگیرید. فرض کنید θ برآورد مبتنی بر x_1, \dots, x_n است. می دانیم که خودگردان ساز ناپارامتری می تواند برای تقریب توزیع $\hat{\theta}$ به کار برده شود. هر تکرار خودگردان ساز بیزی به صورت یک احتمال پسین برای x_i است.

^۶Weng

^۷Mason and Newton

^۸Prestgaard and Wellner

^۹Gasparini

^{۱۰}James

^{۱۱}Banks

^{۱۲}Owen

^{۱۳}Bayesian Bootstrap

لو [۹] ساختار کلی روش BB را به صورت زیر در نظر می‌گیریم، در نهایت $\bar{\mu}^1, \dots, \bar{\mu}^B$ می‌تواند به عنوان تقریب توزیع خودگردان ساز بیزی $\bar{\mu}$ و در نتیجه تقریب توزیع

$(n-1)$ متغیر تصادفی در فاصله $[0, 1]$ در نظر بگیریم

آماره‌های ترتیبی $u_{(1)}, \dots, u_{(n-1)}$ به طوری که $u_{(0)} = 0$

و $u_{(n)} = 1$ را تشکیل می‌دهیم (چرنیک [۴]) و کمیت‌های

$$g_i = u_{(i)} - u_{(i-1)}, \quad i = 1, \dots, n,$$

فراوانی x_i ها در نمونه خودگردان ساز $n_i =$

را محاسبه می‌کنیم. g_i ها اختلاف بین آماره‌های مرتب یکنواخت به صورت بردار $g = (g_1, \dots, g_n)$ هستند که به عنوان بردار احتمالات برای مقادیر x_1, \dots, x_n در هر

تکرار BB در نظر می‌گیریم. برای مثال n مشاهده با نمونه‌گیری با جایگذاری از x_1, \dots, x_n انتخاب شده اند

اما به جای اینکه هر x_i دقیقاً با احتمال $\frac{1}{n}$ انتخاب شوند، از احتمال g_i انتخاب می‌شوند یعنی x_1 از g_1 ، x_2 از g_2 و

... بنابراین با تکرار خودگردان ساز، یک مجموعه جدید از $n-1$ عدد تصادفی یکنواخت از g_i ها داریم. اگر

D_n تابع توزیع تصادفی با وزن‌های g_i برای هر x_i باشد، آن‌گاه با محاسبه $\bar{\theta} = \theta(D_n)$ می‌توان از توزیع شرطی

$R(X, D_n) = \hat{\theta} - \bar{\theta}$ (به شرط X) به عنوان تقریب توزیع پسین $R(\mathbf{X}, F)$ استفاده کرد.

به عنوان مثال فرض کنید θ ، میانگین X است، می‌خواهیم توزیع پسین θ را به دست آوریم. در روش خودگردان ساز

بیزی، بردار $g = (g_1, \dots, g_n)$ را محاسبه می‌کنیم. با فرض اینکه $\theta = \mu$ ، هر تکرار BB یک میانگین به صورت

$\bar{\mu} = \sum_{i=1}^n g_i x_i$ تولید می‌کند. $\bar{\mu} = \sum_{i=1}^n g_i x_i$ روی تمام تکرارهای BB توزیع بوت استرپی میانگین X است که

به عنوان تقریب توزیع پسین μ می‌تواند استفاده شود. حال از روش مونت‌کارلو برای تقریب توزیع BB، $\bar{\mu}$ استفاده می‌کنیم. مرحله بالا را به طور مستقل B بار تکرار

آن‌گاه بردار $n = (n_1, \dots, n_n)$ دارای توزیع چند جمله‌ای است که احتمال $\frac{1}{n}$ را به هر جمله نسبت می‌دهد، و اگر

فراوانی نسبی x_i ها در نمونه خودگردان ساز $f_i =$

آن‌گاه بردار $\mathbf{f} = (f_1, \dots, f_n)$ دارای توزیع چند جمله‌ای $(n, \frac{1}{n}, \dots, \frac{1}{n})$ می‌باشد که

$$E(f_i) = \frac{1}{n},$$

$$Var(f_i) = \frac{n-1}{n^3},$$

$$Cov(f_i, f_j) = -\frac{1}{n^3}, \quad i \neq j,$$

$$Corr(f_i, f_j) = -\frac{1}{n-1}, \quad i \neq j.$$

در روش BB، $g_i = u_{(i)} - u_{(i-1)}$ دارای توزیع $Beta(1, n-1)$ می‌باشد (بهبودیان [۲])، آن‌گاه بردار $g = (g_1, \dots, g_n)$ دارای توزیع دیریکله $(1, \dots, 1)$ است.

(ویلکس [۱۷])

این نمونه است.

۴ خودگردان ساز بیز پارامتری

در حالت پارامتری ارتباط بین خودگردان ساز بیزی و غیر بیزی شفاف تر است. فرآیند خودگردان ساز بیزی پارامتری را براساس یک مساله تک پارامتری ساده مطرح می‌کنیم.

$$E(g_i) = \frac{1}{n},$$

$$Var(g_i) = \frac{n-1}{n^2(n+1)},$$

$$Cov(g_i, g_j) = -\frac{1}{n^2(n+1)}, \quad i \neq j,$$

$$Corr(g_i, g_j) = -\frac{1}{n-1}, \quad i \neq j.$$

اکنون به بررسی مثالی از روبین [۱۵] می‌پردازیم.

مثال ۱.۴. جدول ۱ شامل نمرات ۲۲ دانش‌آموز در دو درس بردار و مکانیک است.

$\hat{\theta} = 0.498$ برآورد ضریب همبستگی نمونه‌ای این دو درس است. می‌خواهیم توزیع پسین پارامتر $\theta_0 = corr(mech, vec)$ را محاسبه کنیم.

فرض می‌کنیم که نمرات هر دانش‌آموز به صورت

$$Y_i = (y_{i1}, y_{i2}) \sim N_2(\mu, \Sigma) \quad ; \quad i = 1, 2, \dots, 22,$$

که $(\hat{\mu}, \hat{\Sigma})$ برآورد درست‌نمایی ماکسیم (MLE) (μ, Σ) هستند. بنابراین از توزیع $N_2(\hat{\mu}, \hat{\Sigma})$ تحت خودگردان ساز پارامتری نمونه‌های تکراری تولید می‌کنیم و برای هر نمونه y^* ، پاسخ خودگردان ساز $\hat{\theta}^*$ را محاسبه می‌کنیم. شکل ۱ بافت نگار ۱۰۰۰۰ تکرار خودگردان ساز برای ضریب همبستگی نمرات دانش‌آموزان را نشان می‌دهد که با تابع توزیع تجربی فیشرف $f_{\theta}(\hat{\theta})$ مقایسه شده است.

$$f_{\theta}(\hat{\theta}) = \frac{(n-2)(1-\theta^2)^{(n-2)/2}(1-\hat{\theta}^2)^{(n-4)/2}}{\pi} \times \int_0^{\infty} \frac{dw}{(\cosh w - \hat{\theta})^{n-1}},$$

مثال ۱.۳. فرض می‌کنیم X متغیری دو حالتی است (0 یا 1)، که پارامتر مورد بررسی (θ) ، بردار احتمال $X = 1$ است، همچنین n_1 را تعداد $x_i = 1$ و $n - n_1$ را تعداد $x_i = 0$ در نظر می‌گیریم.

یک تکرار BB ، ابتدا n احتمال را از n (g_i) تولید شده از $n - 1$ متغیر تصادفی از توزیع یکنواخت در بازه $(0, 1)$ $(U(0, 1))$ به دست آورده و سپس این احتمالات را به مقادیر مشاهده شده X نسبت می‌دهد، بنابراین یک تکرار BB ، n_1 احتمال را به $x_i = 1$

احتمال باقیمانده $(P_1 = \sum_{i=1}^n g_i I(X_i = 1))$ و $n - n_1$ را به $x_i = 0$ نسبت می‌دهد، در نتیجه مقدار BB پارامتر مورد بررسی، P_1 است. تکرار بعدی BB ، به همین روش اما با انتخاب $n - 1$ متغیر جدید از توزیع یکنواخت در بازه $(0, 1)$ $(U(0, 1))$ و محاسبه مقدار جدید P_1 به دست می‌آید، با ادامه تکرارها، توزیع بوت‌استرپی بیزی θ به دست می‌آید.

به دلیل این که g_i ها دارای توزیع $Beta(1, n-1)$ می‌باشند، در نتیجه P_1 که مجموع n_1 ، g_i است دارای توزیع $Beta(n_1, n - n_1)$ است. بنابراین توزیع BB پارامتر مورد بررسی، $Beta(n_1, n - n_1)$ می‌باشد. (روبین [۱۵]) یادآوری می‌کنیم که $Beta(n_1, n - n_1)$ توزیع پسین θ در

جدول ۱: نمرات ۲۲ دانش آموز در دو درس بردار و مکانیک

	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰	۱۱
mech	۷	۴۴	۴۹	۵۹	۳۴	۴۶	۰	۳۲	۴۹	۵۲	۴۴
vec	۵۱	۶۹	۴۱	۷۰	۴۲	۴۰	۴۰	۴۵	۵۷	۶۴	۶۱
	۱۲	۱۳	۱۴	۱۵	۱۶	۱۷	۱۸	۱۹	۲۰	۲۱	۲۲
mech	۳۶	۴۲	۵	۲۲	۱۸	۴۱	۴۸	۳۱	۴۲	۴۶	۶۳
vec	۵۹	۶۰	۳۰	۵۸	۵۱	۶۳	۳۸	۴۲	۶۹	۴۹	۶۳

برآورد MLE ، θ مقدار $\hat{\theta} = 0,498$ است.

عامل تبدیل $R(\theta)$ که نسبت تابع درستنمایی به چگالی خودگردان ساز است را به صورت $R(\theta) = \frac{f_{\hat{\theta}}(\theta)}{f_{\theta}(\theta)}$ تعریف می‌کنیم. در اینجا مقدار $\hat{\theta} = 0,498$ است. تابع احتمال را به صورت

$$P(\theta \in A | \hat{\theta}) = \frac{\int_A \pi(\theta) R(\theta) f_{\hat{\theta}}(\theta) d\theta}{\int_{\Theta} \pi(\theta) R(\theta) f_{\hat{\theta}}(\theta) d\theta}, \quad (1)$$

بازنویسی می‌کنیم. به طور کلی اگر $t(\theta)$ تابعی از θ باشد، امید ریاضی پسین θ به صورت

$$E[t(\theta) | \hat{\theta}] = \frac{\int_{\Theta} t(\theta) \pi(\theta) R(\theta) f_{\hat{\theta}}(\theta) d\theta}{\int_{\Theta} \pi(\theta) R(\theta) f_{\hat{\theta}}(\theta) d\theta}, \quad (2)$$

محاسبه می‌شود.

انتگرال‌های روابط (۱) و (۲) برحسب تابع چگالی خودگردان ساز پارامتری $f_{\hat{\theta}}(\cdot)$ هستند. از آنجایی که

$\theta_1, \dots, \theta_B$ نمونه تصادفی از $f_{\hat{\theta}}(\cdot)$ هستند. انتگرال‌های بالا می‌توانند با میانگین نمونه‌ای برآورد شوند

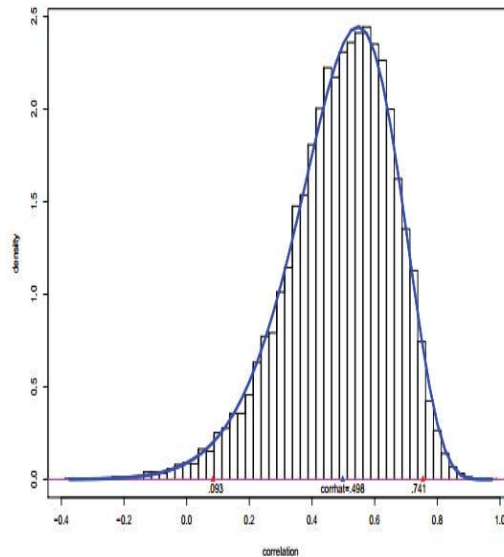
$$\hat{E}[t(\theta) | \hat{\theta}] = \frac{\sum_{i=1}^B t_i \pi_i R_i}{\sum_{i=1}^B \pi_i R_i},$$

اکنون فرض کنید یک چگالی پیشین $\pi(\theta)$ برای پارامتر θ

داریم و می‌خواهیم چگالی پسین $\pi(\theta | \hat{\theta})$ را محاسبه کنیم توسط قانون اعداد بزرگ ثابت می‌شود که $\hat{E}[t(\theta) | \hat{\theta}]$

برای هر زیر مجموعه مانند A از فضای پارامتر $\Theta = \theta$ تقریبی از $E[t(\theta) | \theta]$ است، وقتی که $B \rightarrow \infty$.

در شکل ۲ نمودار پررنگ چگالی پسین $\pi(\hat{\theta} | \theta)$ برای ضریب همبستگی نمرات دانش‌آموزان را نشان می‌دهد.



شکل ۱: مقایسه هیستوگرام با تابع چگالی فیشر

$[-1, 1]$ مطابق قضیه بیز داریم:

$$P(\theta \in A | \hat{\theta}) = \frac{\int_A \pi(\theta) f_{\hat{\theta}}(\theta) d\theta}{\int_{\Theta} \pi(\theta) f_{\hat{\theta}}(\theta) d\theta}.$$

۵ بحث و نتیجه‌گیری

به عقیده روبین خودگردان‌ساز و خودگردان‌ساز بیزی شبیه به هم هستند و ویژگی‌های مشترکی دارند. او معتقد است که محدودیت‌های خودگردان‌ساز بیزی ممکن است ناشی از محدودیت‌های خودگردان‌ساز ناپارامتری باشد. همچنین توصیه می‌کند که اگر چه خودگردان‌ساز بیزی در بعضی از مسایل مفید است اما چون توزیع پیشین (به دلیل ماهیت g_i ها) محدود می‌شود، نباید به عنوان ابزار اصلی استنباط به کار رود. همان‌طور که خودگردان‌ساز ناپارامتری نمی‌تواند به عنوان ابزار اساسی استنباط به کار رود. خودگردان‌ساز پارامتری زمانی که تابع چگالی پارامتر مورد نظر فرم پیچیده‌ای دارد، می‌تواند به عنوان روشی برای برآورد توزیع پسین پارامتر مورد نظر به کار رود.

نمودار خط‌چین چگالی خودگردان‌ساز خام (بدون وزن)، و نمودار مهره دار چگالی خودگردان‌ساز وزن‌دار (BC_a) را نشان می‌دهد، که تقریباً به $\pi(\hat{\theta})$ شبیه است. در توزیع خودگردان‌ساز خام، θ_i در هر تکرار خودگردان‌ساز وزن $\frac{1}{B}$ را می‌گیرد. بنابراین $w_i = \pi_i R_i$ را تعریف می‌کنیم و به این ترتیب برآورد توزیع پسین به صورت

$$P(\theta \in A | \hat{\theta}) = \frac{\sum_{\theta_i \in A} w_i}{\sum_{i=1}^B w_i}, \quad (3)$$

به دست می‌آید.

وزن‌های (BC_a) به صورت $w_i^{BC_a} = \frac{\varphi(z_{\theta_i}/(1+az_{\theta_i})-z_0)}{(1+az_{\theta_i})^2 \varphi(z_{\theta_i}+z_0)}$ که $[z_{\theta_i} = \Phi^{-1} \hat{G}(\theta_i) - z_0]$ تعریف می‌شوند. توزیع پسین به صورت $\pi^{BC_a}(\hat{\theta} | \theta)$ با همان فرمول (۳) برآورد می‌شود.

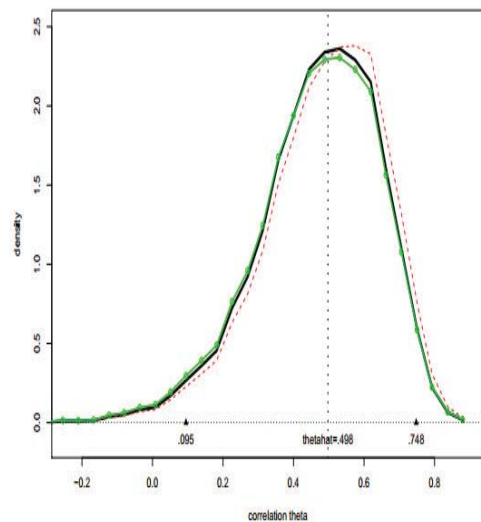
مراجع

[۱] ایران‌پناه، ن. (۱۳۷۷). الگوریتم بوت استرپ بیزی. اندیشه آماری، سال سوم، شماره دوم، ۶۴-۶۹.

[۲] بهبودیان، ج. (۱۳۸۳). آمار ناپارامتری. انتشارات دانشگاه شیراز. چاپ چهارم.

[3] Banks, D. L. (1988). Histospline smoothing the Bayesian bootstrap. *Biometrika*, 75(4), 673-684.

[4] Chernick, M. R. (2007). *Bootstrap Methods: A Guide for Practitioners*



شکل ۲: مقایسه چگالی پسین وزن دار

- [12] Mason, D. M. and Newton, M. A. (1992). A rank statistics approach to the consistency of a general bootstrap. *Annals of Statistics*, 20, 1611-1624.
- [13] Owen, A. B. (1990). Empirical likelihood ratio confidence regions. *Annals of Statistics*, 18, 90-120.
- [14] Praestgaard, D. and Wellner, J. A. (1993). Exchangeably weighted bootstrap of the general empirical process. *Annals of Probability*, 21, 2053-2086.
- [15] Rubin, D. B. (1981). The Bayesian bootstrap. *Annals of Statistics*, 9, 130-134.
- [16] Weng, C. S. (1989). On a second-order asymptotic property of the Bayesian bootstrap mean. *Annals of Statistics*, 17, 705-710.
- [17] Wilks, S. S. (1962). *Mathematical Statistics*, Wiley, New York.
- and Researchers. 2nd ed, Wiley, New York.
- [5] Efron, B. (1979). Bootstrap methods: another look at the jackknife. *Annals of Statistics*, 7, 1-26.
- [6] Efron, B. (2011). The bootstrap and Markov chain Monte Carlo. *Journal of Biopharmaceutical Statistics*, 21(6), 1052-1062.
- [7] Gasparini, M. (1995). Exact multivariate Bayesian bootstrap distribution of moments. *Annals of Statistics*, 23, 762-768.
- [8] James, L. F. (1997). A study of a class of weighted bootstrap for censored data. *Annals of Statistics*, 25, 1595-1621.
- [9] Lo, A. Y. (1987). A large-sample study of the Bayesian bootstrap. *Annals of Statistics*, 15, 360-375.
- [10] Lo, A. Y. (1988). A Bayesian bootstrap for a finite population. *Annals of Statistics*, 16, 1684-1695.
- [11] Lo, A. Y. (1993). A Bayesian bootstrap for censored data. *Annals of Statistics*, 21, 100-123.