

الگوریتم بوت استرپ بیزی

نصرالله ایران‌پناه^۱

چکیده

فرض کنید $X = (X_1, \dots, X_n)$ یک نمونه تصادفی مستقل با توزیعهای یکسان (iid) از توزیع نامعلوم F با پارامتر θ باشد. برآوردکننده‌ی $R(X, F) = \hat{\theta}(X)$ را برای θ و $R(X, F) = \hat{\theta} - \theta$ مانند $R(X, F) = \hat{\theta} - \theta$ در نظر می‌گیریم. افران روش بوت استرپ را با استفاده از طرح نمونه‌گیری مجدد از مشاهدات برای برآورد توزیع نمونه‌گیری $R(X, F)$ و مشخصه‌های آن مانند میانگین و واریانس ارایه کرد.

در استنباط بیزی فرض می‌شود θ یک متغیر تصادفی است که دارای یک توزیع پیشین است و هدف پیدا کردن توزیع پسین $R(X, F)$ است. رویین روش بوت استرپ بیزی ناپارامتری را برای تقریب توزیع پسین $R(X, F)$ ارایه کرده است [۴]. در این مقاله، دو روش بوت استرپ افران و بوت استرپ بیزی رویین ارایه و مقایسه می‌شود. الگوریتم دو روش در یک مثال تشریح و برنامه رایانه‌ای با استفاده از نرم‌افزار آماری S-PLUS ارایه می‌شود.

۱ الگوریتم بوت استرپ افران (BE)

افران در سال ۱۹۷۹ روشی را به نام بوت استرپ برای تقریب توزیع نمونه‌گیری و اندازه‌های دقت (مانند اریبی و واریانس) برآوردکننده‌ها ارایه کرد [۲]. فرض کنید $X = (X_1, \dots, X_n)$ یک نمونه تصادفی (iid) از توزیع نامعلوم F با پارامتر θ باشد. برآوردکننده‌ی $R(X, F) = \hat{\theta}(X)$ را برای θ و $R(X, F) = \hat{\theta} - \theta$ را به صورت تابعی از X و F مانند $R(X, F) = \hat{\theta} - \theta$ در نظر می‌گیریم. هدف، تقریب توزیع نمونه‌گیری $R(X, F)$ و

^۱نصرالله ایران‌پناه، گروه آمار، دانشکده علوم، دانشگاه اصفهان

۱۹۸۱ یک روش بوت استرپ بیزی ناپارامتری را در حالتی که توزیع پیشین نامعلوم است، ارایه کرد. الگوریتم BBR در سه مرحله زیر خلاصه می‌شود:

مرحله ۱ - فرض کنید $U_1, U_2, \dots, U_{n-1}, U_n$ متغیرهای تصادفی (iid) از $(0, 1)$ باشند. آماره‌های مرتب $U_{(0)} = 0 < U_{(1)} < \dots < U_{(n-1)} < U_{(n)} = 1$ و فواصل $(i = 1, \dots, n)$; $g_i = U_{(i)} - U_{(i-1)}$ را در نظر می‌گیریم. اگر D_n تابع توزیع تصادفی با وزنهای g_i برای هر i به صورت زیر باشد:

$$D_n(X) = \sum_{i=1}^n g_i I_{X_i \leq x}$$

آنگاه با محاسبه $\bar{\theta} = \theta(D_n)$ می‌توان از توزیع شرطی $R(X, D_n) = \bar{\theta} - \theta$ (به شرط X) به عنوان تقریب توزیع پسین $R(X, F)$ استفاده کرد. توزیع $R(X, D_n)$ جنبه‌ی نظری دارد و می‌توان آن را از روش مونت‌کارلو تقریب زد.

مرحله ۲ - در روش مونت‌کارلو، مرحله ۱ را به طور مستقل B بار تکرار کرده، نمونه بوت استرپ X^{*b} و از روی آن آماره بوت استرپ $\hat{\theta}^{*b} = \hat{\theta}(X^{*b})$ سپس $R_b = R(X^{*b}, D_n) = \bar{\theta} - \theta^{*b}$ (۳) را محاسبه می‌کنیم.

مرحله ۳ - تابع توزیع تجربی R_1, R_2, \dots, R_B (که جرم احتمال $1/B$ را به هر R_b نسبت می‌دهد)، توزیع $R(X^*, F_n)$ را برای B های بزرگ تقریب می‌زنند.

مشخصه‌های آن، مانند میانگین و واریانس است. الگوریتم در سه مرحله زیر خلاصه می‌شود:

مرحله ۱ - از تابع توزیع تجربی F_n یک نمونه تصادفی (X_1^*, \dots, X_n^*) به صورت (iid) تولید می‌کنیم. به عبارت ساده‌تر نمونه بوت استرپ X^* را به روش نمونه‌گیری تصادفی ساده با جایگذاری از X به دست می‌آوریم. با محاسبه آماره بوت استرپ $\hat{\theta}^* = \hat{\theta}(X^*)$ می‌توان از توزیع $\hat{\theta}^* = \hat{\theta}^* - \bar{\theta}$ به عنوان تقریب توزیع نمونه‌گیری $R(X, F)$ استفاده کرد. توزیع $R(X^*, F_n)$ جنبه‌ی نظری دارد و می‌توان آن را از روش مونت‌کارلو تقریب زد.

مرحله ۲ - در روش مونت‌کارلو، نمونه‌گیری از F_n را به طور مستقل B بار تکرار کرده، نمونه بوت استرپ X^{*b} و از روی آن آماره بوت استرپ $\hat{\theta}^{*b} = \hat{\theta}(X^{*b})$ سپس $R_b = R(X^{*b}, F_n) = \bar{\theta}^{*b} - \bar{\theta}$ (۳) را محاسبه می‌کنیم.

مرحله ۳ - تابع توزیع تجربی R_1, R_2, \dots, R_B (که جرم احتمال $1/B$ را به هر R_b نسبت می‌دهد)، توزیع $R(X^*, F_n)$ را برای B های بزرگ تقریب می‌زنند.

۲ الگوریتم بوت استرپ بیزی روین (BBR)

در استنباط بیزی فرض می‌شود θ یک متغیر تصادفی است که دارای یک توزیع پیشین است و $X = (X_1, \dots, X_n)$ یک نمونه تصادفی (iid) از F (توزیع شرطی X به شرط θ) است. استنباط بیزی بر اساس توزیع پسین $R(X, F)$ یعنی توزیع شرطی $R(X, F)$ به شرط X است. اگر F متعلق به یک خانواده پارامتری با پارامتر θ باشد، توزیع پسین $R(X, F)$ در صورتی قابل محاسبه است که توزیع پیشین معلوم باشد. روین در سال

۳ مقایسه BE و BBR

دو روش BE و BBR را برای پارامتر θ ، میانگین توزیع نامعلوم F یعنی $\theta = \mu = \int x dF(x)$ ، مقایسه می‌کنیم. نمونه تصادفی مشاهده شده x_1, \dots, x_n و برآورد $\hat{\theta} = \bar{x} = \sum x_i/n$ را برای

با به دست آوردن تمام تکرارهای مرحله ۱ (که در عمل کاری غیر ممکن است) و سپس محاسبه $\bar{\mu}$ ، توزیع BBR برای $\bar{\mu}$ به دست می‌آید که به عنوان تقریب توزیع پسین $\bar{\mu}$ می‌تواند استفاده شود. از روش مونت کارلو برای تقریب توزیع بوت استرپ بیزی $\bar{\mu}$ استفاده می‌کنیم. اگر مرحله ۱ الگوریتم BBR را به طور مستقل B بار تکرار کنیم، در نتیجه $\bar{\mu}^B, \dots, \bar{\mu}^1$ به دست می‌آید (مرحله ۲).تابع توزیع $\bar{\mu}^B, \dots, \bar{\mu}^1$ می‌تواند به عنوان تقریب توزیع بوت استرپ بیزی $\bar{\mu}$ و در نتیجه تقریب توزیع پسین $\bar{\mu}$ استفاده شود (مرحله ۳).

الگوریتم دو روش BE و BBR فقط در تخصیص احتمال f_i و g_i به هر x_i باهم اختلاف دارند. در روش BE توزیع نمونه‌گیری برآوردکننده θ شبیه‌سازی و تقریب زده می‌شود، در صورتی که در روش BBR توزیع پسین پارامتر θ شبیه‌سازی و تقریب زده می‌شود.

به سادگی می‌توان نشان داد که دو بردار (f_1, \dots, f_n) و (g_1, \dots, g_n) دارای توزیع \bar{x} جمله‌ای $(\frac{1}{n}, \dots, \frac{1}{n})$ و دیریکله $(1, \dots, 1)$ هستند. و همچنین،

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

$$\text{var}(f_i) = \frac{n-1}{n^2}; \quad \text{var}(g_i) = \frac{n-1}{n^2(n+1)};$$

$$\text{cov}(f_i, f_j) = \text{cov}(g_i, g_j) = \frac{-1}{n-1}.$$

در این بخش دو الگوریتم BE و BBR را در حالتی که یک متغیر دو بعدی $(y, z) = x$ و پارامتر θ ضریب همبستگی جامعه (ρ) است، همراه با یک مثال کاربردی (۱۱، صفحه ۴۲) ارایه می‌کنیم.

۴ یک مثال دو بعدی

از بین دانشجویان رشته الکترونیک دانشگاه صنعتی شریف، برای مطالعه نمرات درک مطلب (y) و دستور زبان (z) درس زبان انگلیسی، نمونه‌ای تصادفی به حجم $15 = n$

در نظر می‌گیریم.

در روش BE، نمونه بوت استرپ x_1^*, \dots, x_n^* را به روش نمونه‌گیری تصادفی ساده با جایگذاری از x_1, \dots, x_n به دست n_i می‌آوریم (مرحله ۱). فرض کنید برای هر $i = 1, 2, \dots, n$ ، $f_i = n_i/n$ فراوانی نسبی x_i باشد که در نمونه بوت استرپ ظاهر شده است و می‌تواند مقادیر $0, 1, \dots, n$ را اختیار کند. واضح است که n_i دارای توزیع $(1, n)$ است. میانگین نمونه بوت استرپ $\bar{x}^* = \sum x_i^*/n = \sum f_i x_i$ را محاسبه می‌کنیم. با به دست آوردن تمام نمونه‌های ممکن بوت استرپ که تعداد آنها $m = n^n$ است و محاسبه میانگین نمونه بوت استرپ روی آنها $\bar{x}^{**}, \dots, \bar{x}^{***}$ به دست می‌آید. توزیع نمونه‌گیری \bar{x} می‌تواند با استفاده از توزیع تجربی $\bar{x}^{**}, \dots, \bar{x}^{***}$ تقریب زده شود. واضح است که به دلیل بزرگ بودن m ، انجام این روش حتی با رایانه نیز کاری مشکل است. در نتیجه از روش مونت کارلو برای تقریب بوت استرپ \bar{x} استفاده می‌کنیم. نمونه‌گیری بوت استرپ در مرحله ۱ را به طور مستقل B بار تکرار کرده و $\bar{x}^{**}, \dots, \bar{x}^{***}$ را محاسبه می‌کنیم (مرحله ۲). تابع توزیع تجربی $\bar{x}^{**}, \dots, \bar{x}^{***}$ می‌تواند به عنوان تقریب توزیع بوت استرپ \bar{x} و در نتیجه تقریب توزیع نمونه‌گیری \bar{x} برای B های بزرگ ($200 \leq B \leq 1000$) استفاده شود (مرحله ۳).

نکته مهم در روش BE این است که در تقریب توزیع

نمونه‌گیری \bar{x} ، فرض شده که تابع توزیع نمونه همان تابع توزیع جامعه است. به عبارت دیگر، نمونه‌های بوت استرپ به طور مستقل از تابع توزیع نمونه (یعنی F_n) تولید شده‌اند.

در روش BBR بردار $(g_1, \dots, g_n) = g$ را مشابه مرحله ۱

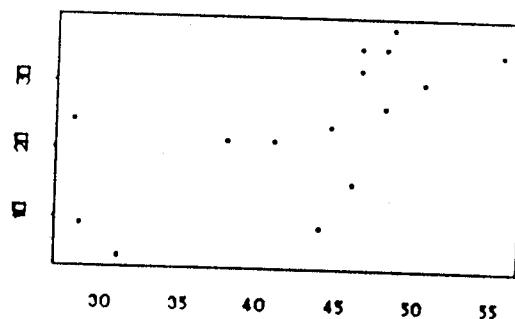
الگوریتم BBR به دست می‌آوریم (مرحله ۱). هر احتمال g_i به صورت یک احتمال پسین برای x_i است. واضح است که g_i دارای توزیع $\text{beta}(1, n - 1)$ است. با فرض $\bar{\mu} = \theta$ ، هر تکرار مرحله ۱، یک میانگین به صورت $\sum g_i x_i = \bar{\mu}$ تولید می‌کند.

همبستگی r را به صورت زیر به دست می‌آوریم:

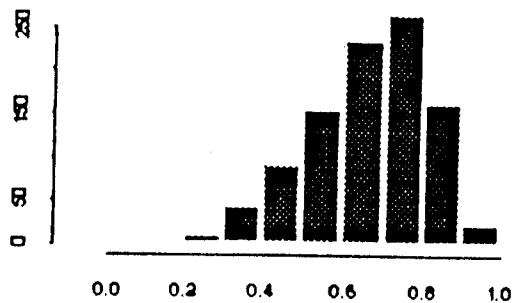
$$\text{bias}^*(r) = \frac{1}{1000} \sum_{b=1}^{1000} r^{*b} - r \\ = 0/662 - 0/663 = 0/001$$

$$SE^*(r) = \frac{1}{999} \sum_{b=1}^{1000} (r^{*b} - \frac{1}{1000} \sum_{b=1}^{1000} r^{*b})^2 \}^{\frac{1}{2}} \\ = 0/152$$

انتخاب و در شکل ۱ نشان داده شده است. ضریب همبستگی نمونه‌گیری $r = 0/663$ است.



شکل ۱



شکل ۲

- همچنین الگوریتم BBR را در این مثال برای تقریب توزیع پسین ρ و مشخصه‌های آن در سه مرحله زیر ارایه می‌کنیم:

مرحله ۱ – بردار تصادفی مشاهده شده x_i را مشابه با مرحله ۱ الگوریتم BBR به دست می‌آوریم.

مرحله ۲ – ضریب همبستگی پسین $\bar{\rho}$ را با استفاده از بردار مرحله ۱، به صورت زیر محاسبه می‌کنیم:

$$\bar{\rho} = \frac{\sum g_i y_i z_i - (\sum g_i y_i)(\sum g_i z_i)}{\{[\sum g_i y_i^2 - (\sum g_i y_i)^2][\sum g_i z_i^2 - (\sum g_i z_i)^2]\}^{\frac{1}{2}}}$$

مرحله ۳ – مراحل ۱ و ۲ را ۱۰۰۰ بار تکرار کرده و ۱۰۰۰ ضریب همبستگی پسین $(\rho^1, \bar{\rho}^1, \dots, \rho^{1000}, \bar{\rho}^{1000})$ را به دست می‌آوریم. شکل ۲ بافت‌نگار توزیع پسین BBR ضریب همبستگی حاصل از ۱۰۰۰

- الگوریتم BE را در این مثال برای تقریب توزیع نمونه‌گیری r و مشخصه‌های آن در سه مرحله زیر ارایه می‌کنیم:

مرحله ۱ – نمونه تصادفی مشاهده شده را به صورت زوجی $15, (y_i, z_i); i = 1, \dots, 15$ در نظر می‌گیریم. از نمونه زوجی (x_1, \dots, x_{15}) به روش نمونه‌گیری تصادفی ساده با جایگذاری، نمونه زوجی بوت استرپ (x_1^*, \dots, x_{15}^*) را به دست می‌آوریم.

مرحله ۲ – ضریب همبستگی نمونه زوجی بوت استرپ r^* را محاسبه می‌کنیم.

مرحله ۳ – مراحل ۱ و ۲ را ۱۰۰۰ بار تکرار کرده و ۱۰۰۰ ضریب همبستگی بوت استرپ $(r^{*1}, \dots, r^{*1000})$ را به دست می‌آوریم. شکل ۲ بافت‌نگار بوت استرپ ضریب همبستگی حاصل از ۱۰۰۰ تکرار مرحله ۳ را نشان می‌دهد که می‌تواند به عنوان تقریب توزیع نمونه‌گیری r استفاده شود. از مشخصه‌های توزیع بوت استرپ استفاده کرده، برآورد بوت استرپ اریبی و خطای معیار ضریب

قابل اجرا است. الگوریتم BBR با استفاده از توابع $gcor$ و $grubin$ قابل اجراست که تابع $gcor$ مقدار ضریب همبستگی پسین ρ را محاسبه می‌کند. در این توابع: x : نمونه تصادفی زوجی، n : حجم نمونه و b : تعداد تکرارهای بوتاسترپ است. خروجی هر دو تابع $bootcor$ و $grubin$ تکرارهای بوتاسترپ r^* و $\bar{\rho}$ است.

تکرار مرحله ۳ را نشان می‌دهد که می‌تواند به عنوان تقریب توزیع نمونه‌گیری پسین ρ استفاده شود. از مشخصه‌های توزیع پسین BBR استفاده کرده، برآورد بوتاسترپ اریبی و خطای معیار ضریب همبستگی پسین ρ را به صورت زیر به دست می‌آوریم:

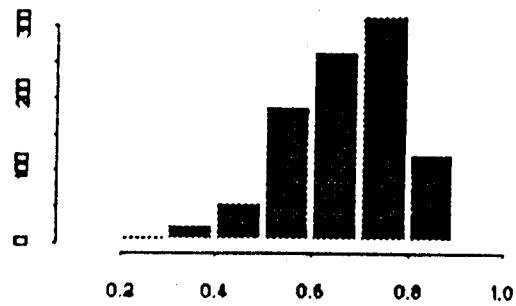
$$bias^*(\rho) = \frac{1}{1000} \sum_{b=1}^{1000} \bar{\rho}^b - r \\ = 0/663 - 0/663 = 0$$

$$SE^*(\rho) = \sqrt{\frac{1}{999} \sum_{b=1}^{1000} (\bar{\rho}^b - \frac{1}{1000} \sum_{b=1}^{1000} \bar{\rho}^b)^2} \\ = 0/121.$$

```
bootcor<-function(n,b,x){
 istar<-matrix(sample(1:n,size=n*b,replace=T)
               ,nrow=b)
  cc<-c(1:b)
  for(i in 1:b){
    xstar<-matrix(c(x[istar[i],1],x[istar[i],2])
                  ,ncal=2)
    cc[i]<-cor(xstar[,1],xstar[,2]))
  }
  return(cc)}
```



```
gcor<-function(g,x){
  y<-sum(g*x[,1])
  z<-sum(g*x[,2])
  yz<-sum(g*x[,1]*x[,2])
  y2<-sum(g*x[,1]*x[,1])
  z2<-sum(g*x[,2]*x[,2])
  rg<-(yz-(y*z))/(sqrt((y2-(y*y)*(z2-(z*z)))))
  return(rg)}
```



شکل ۳

```
grubin<-function(n,b,x){
  rg<-c(1:b)
  g<-c(1:b)
  for(i in 1:b){
    u<-c(0,sort(runif(n-1),1)
    for(j in 1:n){
      g[j]<-(u[j+1]-u[j]))
    rg[i]<-gcor(g,x)}
  return(rg)}
```

۵ برنامه‌های رایانه‌ای

نرم‌افزار S-PLUS یکی از نرم‌افزارهای جدید و توانایی آماری است که امکان برنامه‌نویسی را نیز دارد. مثال بخش ۴ با استفاده از برنامه‌های نوشته شده توسط این نرم‌افزار حل شده است. الگوریتم BE برای پارامتر ρ با استفاده از تابع $bootcor$

مراجع

[۱] ایران پناه، نصرالله و پاشا، عین الله، آشنایی با الگوریتم بوت استرپ، مجله اندیشه آماری، سال دوم، شماره ۱، فروردین ۱۳۷۶.

[۲] Efron, B. (1979), *Bootstrap methods; another look at the Jackknife*, Ann. Statist., 7, 1-26.

[۳] Lo, A., Y. (1987), *A Large sample study of the Bayesian Bootstrap*, Ann. Statist., 15, 360-375.

[۴] Rubin, D. B. (1981), *The Bayesian Bootstrap*, Ann. Statist., 9, 130-134.

[۵] Shao, J. and Tu, D. (1995), *The Jackknife and Bootstrap*, Springer-Verlag.

- هر پدیده‌ای که اتفاق افتاد، مشاهداتی حاصل می‌شود که اگر آنها را ثبت نکنیم، هیچگونه ارزشی ندارند.

«آمار»

- داده‌های جمع‌آوری شده را اگر خلاصه و گویا نکنیم، ارزش آنها ناچیز می‌ماند.

«آمار توصیفی»

- مطالعات انجام شده روی داده‌ها را اگر مورد تجزیه و تحلیل قرار ندهیم و تفسیر نکنیم، ارزش واقعی آنها را از دست می‌دهیم.

«آمار استنباطی»

- اگر نتایج به دست آمده از تجزیه و تحلیل‌های آماری را در برنامه‌ریزی‌ها به کار نبریم، همه‌ی تلاش‌های انجام گرفته بی‌ارزش، و اتلاف وقت بوده است.

«آمار کاربردی»