

محاسبات آماری پیشرفته
ترم اول سال تحصیلی ۹۳
جلسه دهم: روش‌های خودگردان (*Bootstrap Methods*)

حسین باغیشنی

دانشگاه شاهرود

۲ آذر ۱۳۹۳

استنباط آماری با مشاهدات ناقص

آمار، مطالعه روش‌های استخراج استنباط‌هایی از داده‌های ناقص و محدود شده است. مثلاً ما یلیم بدانیم:

- چطور یک نرون در مغز یک موش به کشیده شدن یکی از موهای صورتش، پاسخ می‌دهد

استنباط آماری با مشاهدات ناقص

آمار، مطالعه روش‌های استخراج استنباط‌هایی از داده‌های ناقص و محدود شده است. مثلاً
مایلیم بدانیم:

- چطور یک نرون در مغز یک موش به کشیده شدن یکی از موهای صورتش، پاسخ می‌دهد
- چند تا موش در جنگل ابر زندگی می‌کنند

استنباط آماری با مشاهدات ناقص

آمار، مطالعه روش‌های استخراج استنباط‌هایی از داده‌های ناقص و محدود شده است. مثلاً
مایلیم بدانیم:

- چطور یک نرون در مغز یک موش به کشیده شدن یکی از موهای صورتش، پاسخ می‌دهد
- چند تا موش در جنگل ابر زندگی می‌کنند
- در اردیبهشت، تا چه ارتفاعی آب رودخانه زاینده‌رود زیر پل خواجه بالا خواهد آمد

استنباط آماری با مشاهدات ناقص

آمار، مطالعه روش‌های استخراج استنباط‌هایی از داده‌های ناقص و محدود شده است. مثلاً
مایلیم بدانیم:

- چطور یک نرون در مغز یک موش به کشیده شدن یکی از موهای صورتش، پاسخ می‌دهد
- چند تا موش در جنگل ابر زندگی می‌کنند
- در اردیبهشت، تا چه ارتفاعی آب رودخانه زاینده‌رود زیر پل خواجه بالا خواهد آمد
- متوسط دمای هوا در شاهرود در طول سال، چقدر است
- ...

برای همه این چیزها، مجموعه‌ای از مشاهدات داریم. اما می‌دانیم

برای همه این چیزها، مجموعه‌ای از مشاهدات داریم. اما می‌دانیم

- داده‌ها ناقص هستند
- تکرار آزمایش یا مشاهدات، حتی اگر نهایت سعی خود را برای ثابت نگه داشتن شرایط به کار ببریم، همیشه کم و بیش نتایج متفاوتی خواهد داشت

برای همه این چیزها، مجموعه‌ای از مشاهدات داریم. اما می‌دانیم

- داده‌ها ناقص هستند
 - تکرار آزمایش یا مشاهدات، حتی اگر نهایت سعی خود را برای ثابت نگه داشتن شرایط به کار ببریم، همیشه کم و بیش نتایج متفاوتی خواهد داشت
- بنابراین مضحک است اگر استنباط استخراج شده از داده‌ها را قطعی بدانیم.

مدل‌های تصادفی و تابعی‌ها

اگرچه تکرار یک آزمایش نتایج متفاوتی به دست می‌دهد، بعضی از نتایج نسبت به سایرین بیشتر رخ می‌دهند.

و فراوانی‌های نسبی این پیشامدها، پایدار هستند.

بنابراین، مکانیسم تولید داده‌ها را می‌توان به وسیله توزیع‌های احتمالی و فرآیندهای تصادفی، مدل‌بندی کرد.

مدل‌های تصادفی و تابعی‌ها

اگرچه تکرار یک آزمایش نتایج متفاوتی به دست می‌دهد، بعضی از نتایج نسبت به سایرین بیشتر رخ می‌دهند.

و فراوانی‌های نسبی این پیشامدها، پایدار هستند.

بنابراین، مکانیسم تولید داده‌ها را می‌توان به وسیله توزیع‌های احتمالی و فرآیندهای تصادفی، مدل‌بندی کرد.

کمیت‌هایی مانند مثال‌هایی که ذکر کردیم، به صورت توابعی از مدل تصادفی، یعنی توزیع احتمال زیربنایی، نمایش داده می‌شوند.

چون یک تابع از یک تابع را **تابعی** می‌نامند، و کمیت‌های مورد نظر توابعی از تابع توزیع احتمال واقعی هستند، آن‌ها را تابعی‌های آماری (پارامتر)، $\theta(F)$ ، نیز می‌نامند.

تابعی‌ها می‌توانند عدد حقیقی (مانند تعداد کل موش‌ها)، بردار یا کل یک منحنی (منحنی رگرسیون قد افراد بر روی وزن)، باشند.

استنباط آماری، برآورد کردن این تابعی‌ها یا آزمودن فرضیه‌هایی در مورد آن‌هاست.

مدل‌های تصادفی، عدم قطعیت، استنباط آماری

برآورد پارامترها و سایر استنباط‌ها، توابعی از مقادیر داده هستند.

این به آن معنی است که: **آنها عدم قطعیت همراه با فرآیند تصادفی زیربنایی را به ارث می‌برند.**

اگر آزمایش تکرار شود، ما داده‌های متفاوتی خواهیم داشت اما با یک توزیع قطعی مشخص و اجرای یک روش استنباطی ثابت، نتایج استنباطی متفاوتی در بر خواهد داشت اما دوباره با یک توزیع قطعی مشخص.

مدل‌های تصادفی، عدم قطعیت، استنباط آماری

برآورد پارامترها و سایر استنباط‌ها، توابعی از مقادیر داده هستند.

این به آن معنی است که: **آن‌ها عدم قطعیت همراه با فرآیند تصادفی زیربنایی را به ارث می‌برند.**

اگر آزمایش تکرار شود، ما داده‌های متفاوتی خواهیم داشت اما با یک توزیع قطعی مشخص و اجرای یک روش استنباطی ثابت، نتایج استنباطی متفاوتی در بر خواهد داشت اما دوباره با یک توزیع قطعی مشخص.

آماردانان مایل به استفاده از این توزیع برای کمی‌سازی عدم قطعیت همراه با استنباط‌ها هستند:

مدل‌های تصادفی، عدم قطعیت، استنباط آماری

برآورد پارامترها و سایر استنباطها، توابعی از مقادیر داده هستند.

این به آن معنی است که: **آنها عدم قطعیت همراه با فرآیند تصادفی زیربنایی را به ارث می‌برند.**

اگر آزمایش تکرار شود، ما داده‌های متفاوتی خواهیم داشت اما با یک توزیع قطعی مشخص و اجرای یک روش استنباطی ثابت، نتایج استنباطی متفاوتی در بر خواهد داشت اما دوباره با یک توزیع قطعی مشخص.

آماردانان مایل به استفاده از این توزیع برای کمی‌سازی عدم قطعیت همراه با استنباطها هستند:

- خطای معیار، پاسخی است به این سوال که تا چه حد برآورد پارامتر مورد نظر از یک تکرار آزمایش به دیگری، تغییر می‌کند؟
- یک ناحیه اطمینان برای پارامتر، پاسخ به این سوال است که همه مقادیری از پارامتر که این داده‌ها را، با حداقل یک احتمال مشخص، می‌توانند تولید کنند، کدامند؟

مدل‌های تصادفی، عدم قطعیت، توزیع‌های نمونه‌ای

برای داشتن چیزهایی مانند خطای استاندارد و فاصله اطمینان، نیازمند دانستن توزیع برآوردهای به دست آمده هستیم.

به توزیع برآوردها، توزیع‌های نمونه‌ای گویند.

توزیع‌های نمونه‌ای از توزیع داده‌ها پیروی می‌کنند، زیرا برآوردها توابعی از داده‌ها هستند.

مدل‌های تصادفی، عدم قطعیت، توزیع‌های نمونه‌ای

برای داشتن چیزهایی مانند خطای استاندارد و فاصله اطمینان، نیازمند دانستن توزیع برآوردهای به دست آمده هستیم.

به توزیع برآوردها، توزیع‌های نمونه‌ای گویند.

توزیع‌های نمونه‌ای از توزیع داده‌ها پیروی می‌کنند، زیرا برآوردها توابعی از داده‌ها هستند. از دیدگاه ریاضی، با یک مساله خوش تعریف روبه‌رو هستیم. اما **محاسبه** داستان دیگری است!!!

مدل‌های تصادفی، عدم قطعیت، توزیع‌های نمونه‌ای

برای داشتن چیزهایی مانند خطای استاندارد و فاصله اطمینان، نیازمند دانستن توزیع برآوردهای به دست آمده هستیم.

به توزیع برآوردها، توزیع‌های نمونه‌ای گویند.

توزیع‌های نمونه‌ای از توزیع داده‌ها پیروی می‌کنند، زیرا برآوردها توابعی از داده‌ها هستند.

از دیدگاه ریاضی، با یک مساله خوش تعریف روبه‌رو هستیم. اما محاسبه داستان دیگری است!!!

معمولاً، برآوردها توابع پیچیده‌ای از داده‌ها هستند و صحبت در مورد محاسبه شکل‌های بسته برای توزیع آن‌ها، ناامیدکننده است.

در این موارد، دو راهکار کلاسیک آماردانان عبارتند از:

- تمرکز بر روی موارد خاص و ساده
- استفاده از نظریه مجانبی توزیع‌ها

مدل‌های تصادفی، عدم قطعیت، توزیع‌های نمونه‌ای، محاسبات

تا حدود دهه ۶۰ میلادی، علم آمار در توسعه دو راهکار بالا متمرکز شده بود.

اما با انقلابی که دنیای انفورماتیک و کامپیوترها ایجاد کردند، باعث به کار آمدن مدل‌های آماری پیچیده‌تر (و البته واقعی‌تر) و خارج شدن از دنیای کوچک مدل‌های ساده و سراسر شد.

از طرف دیگر، ممکن است نظریه بزرگ‌نمونه برای مدل‌های پیچیده نیز، درمانی داشته باشد، اما همگرایی به توزیع‌های حدی ممکن است به طور غیرقابل قبولی کند باشد.

تا قبل از دهه ۷۰ میلادی، آمار با مساله کمی‌سازی عدم قطعیت استنباط‌ها بدون در نظر گرفتن پذیره‌های غیرمفید و آمار مجانبی، مواجه شد.

همه راه‌حل‌ها تبدیل شد به محاسبات آماری بیشتر و بیشتر.

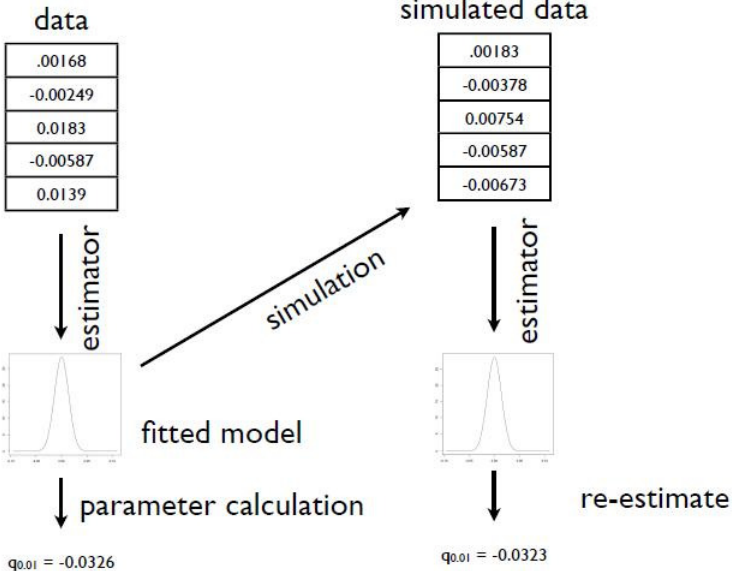
یکی از موفق‌ترین راه‌حل‌ها که توسط بردلی افرون پیشنهاد شد، خودگردان‌سازی، *Bootstrapping*، است.

کلید بررسی عدم قطعیت در برآوردگرها، توزیع نمونه‌ای آن‌ها می‌باشد.

زیرکی افرون در این بود که گفت می‌توانیم تکرار آزمایش را شبیه‌سازی کنیم:

- مدلی را که حدس می‌زنیم مکانیسم واقعی تولید داده‌ها باشد را به داده‌ها برازش می‌دهیم
- اجرای آن مکانیسم، داده‌های شبیه‌سازی شده را تولید می‌کند که دارای توزیع مشابه داده‌های واقعی است
- محاسبه مقدار برآوردگر بر روی داده‌های شبیه‌سازی شده، یک تحقق از توزیع نمونه‌ای برآوردگر را نتیجه می‌دهد
- تکرار این فرآیند به تعداد فراوان، توزیع نمونه‌ای برآوردگر را تقریب می‌زند

چون از مدل برای محاسبه عدم قطعیت خودش استفاده می‌شود، افرون آن را خودگردان نامید.



داده‌های اصلی را با X نشان می‌دهیم. دقت کنید در اینجا منظور از X چارچوب داده است نه یک عدد تنها.

پارامتر برآوردشده به وسیله داده‌ها: $\hat{\theta}$

مجموعه داده‌های شبیه‌سازی شده از مدل برازش شده: $\tilde{X}_1, \dots, \tilde{X}_B$

برآوردهای خودگردان پارامتر، حاصل از داده‌های شبیه‌سازی شده: $\tilde{\theta}_1, \dots, \tilde{\theta}_B$

علاوه بر این فرض کنید، تابعی مورد نظر، $\theta(F)$ ، توسط آماره T به صورت $\hat{t} = T(X)$ برآورد شود و به طور مشابه برای داده‌های شبیه‌سازی شده:

$$\tilde{t}_1 = T(\tilde{X}_1), \dots, \tilde{t}_B = T(\tilde{X}_B),$$

$$\tilde{t} = (\tilde{t}_1, \dots, \tilde{t}_B).$$

آماره T ممکن است تابعی مستقیم از پارامترهای برآوردشده باشد، یا تابعی غیرمستقیم.

در این جا فرض می‌کنیم مدل برای مقداری از θ یک مدل درست است. این مقدار را با θ نشان می‌دهیم. به طور مشابه فرض کنید مقدار درست تابعی، t ، است.

واریانس و خطای استاندارد

ساده‌ترین معیار برای کمی‌سازی عدم قطعیت، واریانس یا خطای استاندارد است:

$$\begin{aligned}\hat{Var}(\hat{t}) &= Var(\tilde{t}) \\ \hat{se}(\hat{t}) &= sd(\tilde{t})\end{aligned}$$

منطق این روابط اینست که \tilde{X} شبیه‌سازی شده توزیع یکسانی با متغیر واقعی X دارد که داده‌های اصلی، x ، از آن استخراج شده‌اند. بنابراین اجرای مشابه روش برآورد بر روی داده‌های شبیه‌سازی شده، توزیع نمونه‌ای را نتیجه خواهد داد.

دقت کنید در این جا فرض بر آن است که مدل در نظر گرفته شده درست است و $\hat{\theta}$ خیلی از مقدار واقعی θ دور نیست.

```
rboot <- function(B, statistic, simulator, ...) {  
  tboots <- replicate(B, statistic(simulator(...)))  
  return(tboots)  
}  
  
bootstrap.se <- function(simulator, statistic, B, ...) {  
  tboots <- rboot(B, statistic, simulator, ...)  
  se <- sd(tboots)  
  return(se)  
}
```

تصحیح اریبی

می‌توان از روش خودگردان برای تصحیح اریبی یک برآوردگر اریب، استفاده کرد. از آنجا که توزیع نمونه‌ای \tilde{t} نزدیک به توزیع \hat{t} است، و \hat{t} نیز نزدیک به t ، پس

$$\mathbb{E}(\hat{t}) - t \approx \mathbb{E}(\tilde{t}) - \hat{t}.$$

سمت چپ تساوی، اریبی هست که به دنبال محاسبه آن هستیم و سمت راست مقداری است که از روی نمونه‌های خودگردان قابل دستیابی است.

دقت کنید که تساوی بالا مادامی معتبر است که توزیع نمونه‌ای $t - \hat{t}$ نزدیک به توزیع نمونه‌ای $\hat{t} - \tilde{t}$ باشد. این لزوم، شرط ضعیفتری نسبت به نزدیک بودن توزیع‌های نمونه‌ای \tilde{t} به \hat{t} و \hat{t} به t است.

یک شرط کافی (اما نه لازم) برای برقراری تساوی بالا، آن است که کمیت $t - \hat{t}$ محوری یا تقریباً محوری باشد.

```
bootstrap.bias <- function(simulator, statistic, B, t.hat, ...) {  
  tboots <- rboot(B, statistic, simulator, ...)  
  bias <- mean(tboots) - t.hat  
  return(bias)  
}
```

$$Pr(t. \in C) = 1 - \alpha,$$

در زمان محاسبه فاصله اطمینان، عدم شناخت دقیق توزیع نمونه‌ای به آن معنی است که سطح اطمینان واقعی (نرخ پوشش فاصله) مقدار دلخواه $1 - \alpha$ نخواهد بود.

هر چه نرخ پوشش تجربی نزدیک‌تر به مقدار اسمی باشد، به معنی تقریب بهتر و در نتیجه فاصله اطمینان دقیق‌تر است.

فرض کنید C_l و C_u به ترتیب کران‌های پایین و بالای فاصله اطمینان باشند. برای فواصل اطمینان با دم‌های برابر داریم:

$$\begin{aligned} \frac{\alpha}{2} &= Pr(C_l \geq t_0) = Pr(C_l - \hat{t} \geq t_0 - \hat{t}) \\ &= Pr(\hat{t} - C_l \leq \hat{t} - t_0) \\ \frac{\alpha}{2} &= Pr(\hat{t} - C_u \geq \hat{t} - t_0) \end{aligned}$$

فاصله اطمینان خودگردانی پایه

روش خودگردانی، توزیع $\tilde{t} - \hat{t}$ را نتیجه می‌دهد که تقریباً با توزیع $t - \hat{t}$ یکی است. با محاسبه این توزیع و \hat{t} ، می‌توان کران‌های C_l و C_u را به دست آورد:

$$C_l = \hat{t} - \left(Q_{\tilde{t}}\left(1 - \frac{\alpha}{2}\right) - \hat{t} \right)$$

$$C_u = \hat{t} - \left(Q_{\tilde{t}}\left(\frac{\alpha}{2}\right) - \hat{t} \right)$$

که در آن $Q_{\tilde{t}}$ تابع چنک \tilde{t} است. یعنی چنک‌های توزیع نمونه‌ای نمونه‌های شبیه‌سازی شده را به دست می‌دهد.

تمرین: دو رابطه بالا را اثبات کنید.

این فاصله، یک فاصله اطمینان خودگردانی پایه یا فاصله اطمینان مبتنی بر کمیت محوری است. به سادگی قابل محاسبه است و دقت آن هم قابل قبول است.

```
bootstrap.ci.basic <- function(simulator, statistic, B, t.hat,  
                               alpha, ...) {  
  tboots <- rboot(B,statistic, simulator, ...)  
  ci.lower <- 2*t.hat - quantile(tboots,1-alpha/2)  
  ci.upper <- 2*t.hat - quantile(tboots,alpha/2)  
  return(list(ci.lower=ci.lower,ci.upper=ci.upper))  
}
```


فواصل اطمینان خودگردانی استودنت شده

فاصله اطمینان خودگردانی پایه، مبتنی بر پذیره تقریبا یکسان بودن توزیع $\hat{t} - t$ با توزیع $\tilde{t} - t$ است.

اما حتی زمانی که این پذیره نادرست باشد، توزیع

$$\tau = \frac{\hat{t} - t}{\hat{se}(\hat{t})}$$

نزدیک به توزیع

$$\tilde{\tau} = \frac{\tilde{t} - t}{se(\tilde{t})}$$

خواهد بود. این کمیت‌ها شبیه آماره t در آزمون t -استودنت هستند و چون آزمون t توسط استودنت کشف شد، به آن‌ها کمیت‌های استودنت شده می‌گویند.

فواصل اطمینان خودگردانی استودنت شده: ادامه

اگر τ و $\tilde{\tau}$ توزیع یکسانی داشته باشند، می توان فاصله اطمینان زیر را به دست آورد:

$$\left(\hat{t} - \hat{se}(\hat{t}) Q_{\tilde{\tau}} \left(1 - \frac{\alpha}{2}\right), \hat{t} + \hat{se}(\hat{t}) Q_{\tilde{\tau}} \left(\frac{\alpha}{2}\right) \right)$$

این فاصله با فاصله پایه یکی خواهد بود هرگاه $se(\tilde{t}) = \hat{se}(\hat{t})$ و در غیر این صورت متفاوت خواهند بود.

برای محاسبه $se(\tilde{t})$ نیاز به سطح دومی از خودگردان سازی داریم. به الگوریتم زیر دقت کنید:

فواصل اطمینان خودگردانی استودنت شده: الگوریتم

۱ مدل را با $\hat{\theta}$ برازش بده و \hat{t} را محاسبه کن

۲ برای $i = 1, \dots, B_1$:

الف) \tilde{X}_i را از $\hat{\theta}$ تولید کن

ب) مقادیر $\tilde{\theta}_i$ و \tilde{t}_i را محاسبه کن

ج) برای $j = 1, \dots, B_2$:

۱. X_{ij}^\dagger از $\tilde{\theta}_i$ تولید کن

۲. t_{ij}^\dagger را محاسبه کن

د) $\tilde{\sigma}_i$ را برابر انحراف استاندارد t_{ij}^\dagger ها قرار بده

ه) برای تمام j ها قرار بده $\tilde{\tau}_{ij} = \frac{t_{ij}^\dagger - \tilde{t}_i}{\tilde{\sigma}_i}$

۳ $\hat{se}(\hat{t})$ را برابر انحراف استاندارد \tilde{t}_i ها قرار بده

۴ چندک های $\alpha/2$ و $1 - \alpha/2$ توزیع $\tilde{\tau}$ را محاسبه کن

۵ همه کمیت های مورد نظر را در فاصله اطمینان جایگذاری کن

فواصل اطمینان خودگردانی چندکی

مزیت فواصل استودنت شده نسبت به فواصل پایه، دقت بالاترشان است.
عیب آن‌ها نیز زمان بر بودن محاسبات است.

نوع دیگری از فواصل اطمینان خودگردانی، فواصل چندکی است که به راحتی به صورت زیر تشکیل می‌شود:

$$(Q_{\bar{t}}(\alpha/2), Q_{\bar{t}}(1 - \alpha/2)).$$

این نوع فواصل اطمینان به سادگی قابل محاسبه‌اند، اما دقیق نیستند.

البته همه این فواصل اطمینان مطرح شده، صورت‌های مختلفی دارند که توسط افراد مختلف پیشنهاد شده‌اند.

آزمون فرضیه خودگردانی

برای آزمون فرضیه‌ها، دو توزیع نمونه‌ای متفاوت را محاسبه می‌کنیم:

- توزیع آماره آزمون تحت فرضیه صفر که برای محاسبه اندازه آزمون و سطح معنی‌داری است

- توزیع آماره آزمون تحت فرضیه جانشین که برای محاسبه توان آزمون است

هر دو توزیع با روش خودگردان‌سازی قابل محاسبه هستند.

در آزمون فرضیه، آماره مورد نظر، t ، آماره آزمون می‌باشد.

```
boot.pvalue <- function(test,simulator,B,testthat, ...) {  
  testboot <- rboot(B=B, statistic=test, simulator=simulator, ...)  
  p <- (sum(test >= testthat)+1)/(B+1)  
  return(p)  
}
```

این کد، p - مقدار خودگردانی را برای یک آزمون محاسبه می‌کند. دقت کنید که $testthat$ مقدار آماره آزمون برای داده‌های اصلی است و $test$ تابعی است که معرف آماره آزمون است.

برای محاسبه توان آزمون به روش مشابه می‌توان عمل کرد. فقط باید نمونه‌های خودگردانی تحت فرضیه جانشین تولید شوند و $testthat$ مقدار بحرانی آزمون باشد نه مقدار مشاهده‌شده آماره آزمون.

مثال: قانون نابرابری ثروت پارتو

توزیع پارتو (توزیع قانون توانی) یک مدل معمول برای داده‌های با دم‌های سنگین است. به این معنی که چگالی احتمال $f(x)$ زمانی که $x \rightarrow \infty$ خیلی کند به صفر میل می‌کند. یا به عبارت دیگر، این توزیع به شدت به راست چوله است و میانگین آن خیلی بزرگتر از میانه است.

$$f(x) = \frac{\theta - 1}{x} \left(\frac{x}{x_0} \right)^{-\theta}$$

که در آن x_0 مقیاس می‌نیم توزیع است.

تمرین. نشان دهید که x_0 مد توزیع پارتو است.

اگر x_0 معلوم باشد، آن گاه:

$$\hat{\theta} = 1 + \frac{n}{\sum_{i=1}^n \log \frac{x_i}{x_0}}$$

که برآوردگری سازگار و کاراست.

مثال پارتو: ادامه

فایل *pareto.R* شامل تعدادی تابع مرتبط با توزیع پارتو است. یکی از توابع آن *pareto.fit* است که مدل بالا را به داده‌ها برازش می‌دهد.

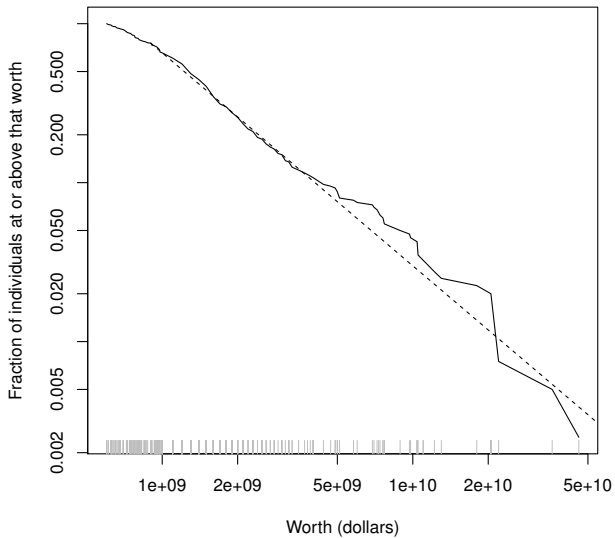
پارتو، وقتی که می‌خواست توزیع داده‌های ثروت را مدل‌بندی کند، به چگالی پارتو رسید.

تقریباً در کشورهای مختلف و زمان‌های متفاوت، دم بالایی توزیع درآمد و ثروت از یک قانون توانی پیروی می‌کند.

پارامتر θ در این توزیع با مقدار پولی که کم و بیش در دست ثروتمندترین مردم کشور هست، تغییر می‌کند.

شکل صفحه بعد، توزیع ثروت را برای ۴۰۰ پولدار اول آمریکا در سال ۲۰۰۳ نشان می‌دهد. با در نظر گرفتن $x_i = 9 \times 10^8$ ، تعداد افرادی که در دم قرار می‌گیرند، ۳۰۲ نفر است و $\hat{\theta} = 2/34$:

```
source("pareto.R")
wealth <- scan("wealth.dat")
wealth.pareto <- pareto.fit(wealth,threshold=9e8)
> signif(wealth.pareto$exponent,3)
[1] 2.34
```



چه میزان عدم قطعیت در این برآورد θ وجود دارد؟

با روش خودگردانی می‌توان به این سوال پاسخ داد:

```
rboot.pareto <- function(B,exponent,x0,n) {
  replicate(B,pareto.fit(rpareto(n,x0,exponent),x0)$exponent)
}
pareto.se <- function(B,exponent,x0,n) {
  return(sd(rboot.pareto(B,exponent,x0,n)))
}
pareto.bias <- function(B,exponent,x0,n) {
  return(mean(rboot.pareto(B,exponent,x0,n)) - exponent)
}
```

با $\hat{\theta} = 2/34$ ، $x_0 = 9 \times 10^8$ ، $n = 302$ و $B = 10^4$ ، خطای معیار (عدم قطعیت برآورد) برابر $0/077$ به دست می‌آید که با نتایج مجانبی هم هم‌خوانی دارد.

به طور مجانبی نیز (با توجه به سازگاری برآوردگر) اریبی به سمت صفر میل می‌کند. میزان اریبی خودگردان نیز برابر $10^{-3} \times 4$ به دست آمده است که قابل صرف‌نظر است.

می توان فاصله اطمینان خودگردانی پایه را نیز محاسبه کرد:

```
pareto.ci <- function(B,exponent,x0,n,alpha) {  
  tboot <- rboot.pareto(B,exponent,x0,n)  
  ci.lower <- 2*exponent - quantile(tboot,1-alpha/2)  
  ci.upper <- 2*exponent - quantile(tboot,alpha/2)  
  return(list(ci.lower=ci.lower, ci.upper=ci.upper))  
}
```

با استفاده از همان مشخصات قبلی، فاصله ۹۵٪ پایه به صورت $(2/18, 2/48)$ به دست می آید.

روش خودگردان سازی ناپارامتری

خودگردان سازی، توزیع نمونه‌ای را همراه با سه منبع خطای برآورد، تقریب می‌زند:

- ۱ خطای شبیه‌سازی: استفاده از تعداد تکرار فراوان اما متناهی برای دستیابی به توزیع نمونه‌ای کامل. تعداد کافی تکرار و طراحی مناسب شبیه‌سازی، می‌تواند این خطا را به اندازه دلخواه کوچک کند.
- ۲ خطای آماری: توزیع نمونه‌ای پارامترهای برآورد شده خودگردانی تحت مدل برازش شده، دقیقاً با توزیع نمونه‌ای پارامترهای برآورد شده تحت مدل واقعی مکانیسم تولید داده‌ها یکی نیست. توزیع نمونه‌ای با تغییر پارامترها، تغییر می‌کند و برآورد اولیه ما کاملاً دقیق نیست. اما اغلب مشخص شده است که توزیع برآوردها حول مقدار واقعی پایاتر از توزیع خود برآوردهاست. بنابراین کم کردن برآورد اولیه از مقادیر خودگردان، باعث کاهش خطای آماری می‌شود.
- ۳ خطای مشخص سازی: داده‌ها از مدل‌هایی که ما در نظر می‌گیریم دقیقاً پیروی نمی‌کنند. بنابراین شبیه‌سازی مدل هرگز با توزیع نمونه‌ای واقعی یکی نمی‌شود.

خطای آماری:

$$F_n \longrightarrow \tilde{X} \longrightarrow F_n^*,$$

$$F_n \longrightarrow \tilde{t} \longrightarrow F_n^*$$

خطای مشخص سازی:

$$F \longrightarrow X \longrightarrow F_n$$

$$F \longrightarrow t \longrightarrow F_n$$

برای درک خطای آماری، به مثال ۷.۱ کتاب توجه کنید.

روش خودگردان‌سازی ناپارامتری: ادامه

افرون، یک ایده معرکه دومی هم داشت، که مربوط به خطای مشخص‌سازی می‌شد.

وی پیشنهاد کرد برای فرار از این خطا، به جای شبیه‌سازی از مدل، از خود داده‌ها (یا به عبارتی توزیع تجربی داده‌ها) باز نمونه‌گیری کنیم.

در واقع خود داده‌ها (یا توزیع تجربی داده‌ها) آماره بسنده برای پارامتر است.

از دیدگاه دیگر، توزیع تجربی داده‌ها کم‌تعصب‌ترین برآورد ممکن از توزیع واقعی زیربنایی است. هر چیز دیگری ارزیابی و پیش‌داوری وارد می‌کند. ممکن است مدل انتخاب‌شده دقیق باشد، اما احتمال نامناسب بودن آن نیز کم نیست.

بسیاری از کمیت‌ها را می‌توان به طور مستقیم، و بدون وساطت یک مدل پارامتری، از روی توزیع تجربی برآورد کرد.

خودگردان‌سازی ناپارامتری افرون، مجموعه داده اصلی را به عنوان یک جامعه متناهی در نظر می‌گیرد و یک نمونه شبیه‌سازی‌شده جدید، که در آن هر مشاهده با احتمال برابر انتخاب می‌شود، از آن استخراج می‌کند. سپس برآورد را با نمونه جدید تکرار می‌کند.

در واقع، معمولاً وقتی آماردانان اسم خودگردان را می‌آورند، منظورشان این نوع خودگردان است.

الگوریتم خودگردان سازی ناپارامتری

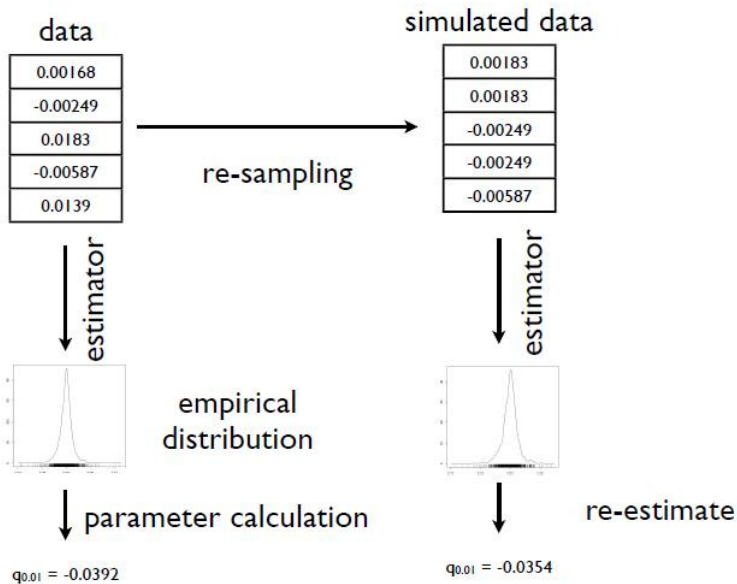
توزیع نمونه‌ای \hat{t} را می‌توان به کمک روش خودگردان سازی به صورت زیر تقریب زد:

۱ برای هر تکرار خودگردان $b = 1, \dots, B$:

الف) نمونه $\tilde{x}^{(b)} = (\tilde{x}_1, \dots, \tilde{x}_n)$ را با نمونه‌گیری با جایگذاری از نمونه مشاهده شده x_1, \dots, x_n تولید کن

ب) مقادیر $\tilde{t}^{(b)}$ و $\tilde{\theta}^{(b)}$ را محاسبه کن

۲ برآورد توزیع نمونه‌ای خودگردان، F_n^* ، توزیع تجربی $\tilde{t}^{(1)}, \dots, \tilde{t}^{(B)}$ می‌باشد



روش خودگردان سازی ناپارامتری: ادامه

هر آنچه که با خودگردان سازی پارامتری مطرح شد، با خودگردان سازی ناپارامتری نیز قابل طرح و انجام است.

تنها تفاوت در توزیع داده های شبیه سازی شده خودگردانی است که همان توزیع تجربی داده های اصلی است.

```
resample <- function(x) { sample(x,size=length(x),replace=TRUE) }
resamp.pareto <- function(B,data,x0) {
  replicate(B,pareto.fit(resample(data),threshold=x0)$exponent)
}
> sd(resamp.pareto(B,wealth,x0))
[1] 0.07762305
> mean(resamp.pareto(B,wealth,x0))-exponent
[1] 0.003534574
resamp.pareto.CI <- function(B,data,alpha,x0) {
  thetahat <- pareto.fit(data,threshold=x0)$exponent
  thetaboot <- resamp.pareto(B,data,x0)
  ci.lower <- thetahat - (quantile(thetaboot,1-alpha/2) - thetahat)
  ci.upper <- thetahat - (quantile(thetaboot,alpha/2) - thetahat)
  return(list(ci.lower=ci.lower,ci.upper=ci.upper))
}
> resamp.pareto.CI(B,wealth,0.05,x0)
$ci.lower
 97.5%
2.173528

$ci.upper
  2.5%
2.476391
```


خودگردان سازی پارامتری در مقابل ناپارامتری

اگر یک مدل مشخص (پارامتری) مناسب داشته باشیم، شبیه سازی از آن مدل نتایج دقیق تری را، با حجم نمونه ثابت n ، نسبت به بازنمونه گیری از توزیع تجربی داده ها (ناپارامتری) به دست خواهد داد.

در واقع، در این حالت، برآورد پارامتری توزیع سریع تر از توزیع تجربی، به توزیع واقعی همگرا می شود.

در مقابل اگر مدل پارامتری به اشتباه مشخص شده باشد، به سرعت به یک توزیع نادرست همگرا می شود. در چنین مواردی استفاده از روش ناپارامتری، نتایج خیلی بهتری به دست خواهد داد.

معمولا، چون در اغلب کاربردها نسبت به پذیره های مدل های پارامتری تردید زیادی وجود دارد، ترجیح داده می شود از روش های بازنمونه گیری استفاده شود. مگر آن که بتوان قانع شد که یک مدل پارامتری، تقریب خیلی خوبی از واقعیت پدیده تصادفی خواهد بود.

مثال: مجموعه داده آزمون ورودی دانشکده حقوق

در بسته *bootstrap* در *R*، که مربوط به مباحث و مثال‌های کتاب افرون و تیشیرانی است، مجموعه داده‌ای به نام داده‌های دانشکده حقوق وجود دارد که شامل معدل نمره آزمون‌های ورودی ۸۲ دانشکده حقوق، *LSAT*، و معدل نمرات دبیرستان شرکت‌کنندگان در آزمون، *GPA*، است.

این داده‌ها در این بسته با نام *law82* قابل دسترسی است. یک نمونه ۱۵ تایی از این داده‌ها در صفحه ۱۸۵ کتاب آورده شده است.

هدف برآورد ضریب وابستگی بین این دو متغیر و محاسبه برآورد خودگردانی خطای استاندارد این برآورد می‌باشد.

در این مثال تابعی *t*، ضریب همبستگی بین دو متغیر *LSAT* و *GPA* می‌باشد.

```
library(bootstrap) # for the law data
> print(cor(law$LSAT, law$GPA)) # for a sample with size 15
[1] 0.7763745
> print(cor(law82$LSAT, law82$GPA))
[1] 0.7599979
```

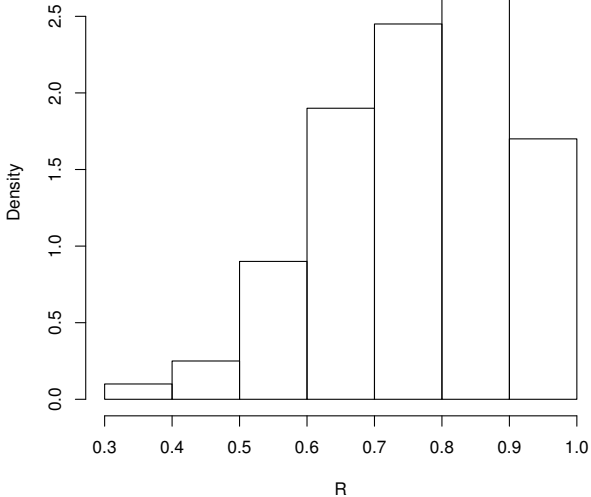
ادامه مثال: مجموعه داده آزمون ورودی دانشکده حقوق

```
B <- 200 # number of replicates
n <- nrow(law) # sample size
R <- numeric(B) # storage for replicates
# bootstrap estimate of standard error of R
for (b in 1:B) {
  # randomly select the indices
  i <- sample(1:n, size = n, replace = TRUE)
  LSAT <- law$LSAT[i]      #i is a vector of indices
  GPA <- law$GPA[i]
  R[b] <- cor(LSAT, GPA)
}
# output
> print(se.R <- sd(R))
[1] 0.1348793
```

برآورد خطای استاندارد بر اساس نظریه نرمال برابر 0.115 است. (چرا؟)

تمرین: با تابع *boot* در بسته *boot* می‌توان محاسبات بالا را انجام داد. برای مثال آزمون ورودی دانشکده حقوق، محاسبات را تکرار کنید و با نتایج بالا مقایسه کنید.

Histogram of R



ادامه مثال: مجموعه داده آزمون ورودی دانشکده حقوق

برای برآورد اریبی نیز می‌توان مشابه مثال‌های قبلی عمل کرد:

```
# sample estimate for n=15
theta.hat <- cor(law$LSAT, law$GPA)
# bootstrap estimate of bias
B <- 2000 # larger for estimating bias
n <- nrow(law)
theta.b <- numeric(B)
for (b in 1:B) {
  i <- sample(1:n, size = n, replace = TRUE)
  LSAT <- law$LSAT[i]
  GPA <- law$GPA[i]
  theta.b[b] <- cor(LSAT, GPA)
}
bias <- mean(theta.b - theta.hat)
> bias
[1] -0.002769211
```

برآورد اریبی برای یک برآوردگر نسبت

داده‌های *patch* نیز در بسته *bootstrap* قرار دارد که توضیح کامل آن در کتاب افرون و تیشیرانی (۱۹۹۳) هست.

این داده‌ها شامل اندازه‌های یک هرمون خاص در گردش خون ۸ آزمودنی بعد از پانسما ن پزشکی است. پارامتر مورد نظر عبارتست از:

$$\theta = \frac{E(new) - E(old)}{E(old) - E(placebo)}$$

که در آن $E(new)$ و $E(old)$ پانسما ن‌های قدیمی و جدید هستند. برای این داده‌ها، اگر $|t| \leq 0.20$ ، نشان‌دهنده برابری عملکرد پزشکی پانسما ن‌های قدیم و جدید است.

تابعی مورد نظر آماره $t = \frac{\bar{Y}}{Z}$ است. هدف برآورد خودگردانی اریبی آماره نسبت برابری عملکرد است.

```
data(patch, package = "bootstrap")
> patch
```

	subject	placebo	oldpatch	newpatch	z	y
1	1	9243	17649	16449	8406	-1200
2	2	9671	12013	14614	2342	2601
3	3	11792	19979	17274	8187	-2705
4	4	13357	21816	23798	8459	1982
5	5	9055	13850	12560	4795	-1290
6	6	6290	9806	10157	3516	351
7	7	12412	17208	16570	4796	-638
8	8	18806	29044	26325	10238	-2719

برآورد اریبی برای یک برآوردگر نسبت: ادامه

```
n <- nrow(patch) # in bootstrap package
B <- 2000
theta.b <- numeric(B)
theta.hat <- mean(patch$y) / mean(patch$z)
# bootstrap
for (b in 1:B) {
  i <- sample(1:n, size = n, replace = TRUE)
  y <- patch$y[i]
  z <- patch$z[i]
  theta.b[b] <- mean(y) / mean(z)
}
bias <- mean(theta.b) - theta.hat
se <- sd(theta.b)
> print(list(est=theta.hat, bias = bias, se = se, cv = bias/se))
$est
[1] -0.0713061

$bias
[1] 0.004647646

$se
[1] 0.09888232

$cv
[1] 0.04700179
```

معمولا اگر $\frac{|bias|}{se} \leq 0.25$ ، نیازی به تصحیح اریبی نیست. در این مثال این نسبت کمتر از 0.05 است.