

«پروژه درس استنباط یک»

روش مونت کارلو با زنجیره‌ی مارکف

MCMC

استاد : دکتر بهالدین خالدی

محی الدین ایزدی عبدالله حسنی جلیلیان

۱۳۸۴ بهمن

فصل ۱

روش مونت کارلو

۱.۱ مقدمه

تعريف

روش مونت کارلو راه حل تقریبی مسائل ریاضی ، آمار ، فیزیک و... است که برای شبیه سازی کمیت های تصادفی برای تقریب (به عبارت بهتر برآورد) کمیتی مجهول به کار می رود. در واقع هر روشی که سعی بر حل مسئله با تولید اعداد تصادفی دارد را روش مونت کارلو گویند.

تاریخچه پیدایش روش مونت کارلو

این روش در سال 1949 با انتشار مقاله‌ای تحت عنوان روش مونت کارلو از دو ریاضیدان آمریکایی به نام های J.NEYMAN و S.ULAM شناخته شد . پایه های نظری این روش از مدت ها قبل مشخص شده بود و در قرن نوزدهم و اوایل قرن بیستم گاهی اوقات مسائل آماری را با کمک شبیه سازی کمیت های تصادفی که همان روش مونت کارلو است حل می کردند.اما به طور قطع تا پیش از اختراع رایانه های الکترونیکی از این روش به دلیل خسته کننده بودن شبیه سازی کمیت های تصادفی زیاد استفاده نمی شد . بنابراین می توان شروع استفاده از روش مونت کارلو به عنوان یک تکنیک عددی عمومی را همزمان با اختراع رایانه های الکترونیکی دانست . نام مونت کارلو از شهر مونت کارلو ، در موناکو که به خاطر قمارخانه هاییش شهرت دارد گرفته شده است . مسائلی که می توان آن ها را به روش مونت کارلو حل کرد :

روش مونت کارلو ، شبیه سازی هر روندی که توسط عامل های تصادفی ایجاد شده اند را ممکن می سازد. اما این تنها مورد استفاده روش مونت کارلو نیست. برای بیشتر مسائل ریاضی که در آن ها شناس دخالت ندارد ما می توانیم به طور مصنوعی مدل احتمال آن را با تعداد زیادی نمونه تعیین کنیم ، به این دلیل روش مونت کارلو را یک روش عمومی حل مسائل ریاضی می توان دانست .

۲.۱ محاسبه انتگرال ها با روش مونت کارلو

فرض کنیم g یک تابع حقیقی مقدار باشد که علاقمند به محاسبه انتگرال g روی $[a, b]$ هستیم .

$$\theta = \int_a^b g(x)dx < \infty$$

برای محاسبه این انتگرال می توان به چند روش عمل کرد:

۱) روش تحلیلی : یافتن پادمشتق g و با استفاده از قضیه اساسی حساب دیفرانسیل و انتگرال ،

$$\theta = \int_a^b g(x)dx = G(b) - G(a)$$

۲) روش های عددی : اگر قادر نباشیم مسئله ها را به روش تحلیلی حل کنیم از روش های عددی که برای تقریب مقدار θ وجود دارند ، استفاده می کنیم . در این روش عمدتا تابع g با توابع ساده تری مانند توابع خطی تکه ای ، چند جمله ای و ...، تقریب زده شده و مقدار انتگرال برای این توابع خوش رفتار محاسبه شده ، به عنوان تقریب θ ارائه می شود.

۳) روش مونت کارلو : این روش نیز یک روش عددی است که در مواردی که روش های عددی در قسمت بالا با پیچیدگی های زیادی همراه هستند یا خطای این روش زیاد است و یا پیش فرض های روش های عددی برقرار نیستند ، به ویژه وقتی که g یک بعدی نباشد کاربرد دارد.

۱.۲.۱ روش (براورد) مونت کارلو:

هدف محاسبه انتگرال زیر است

$$\theta = \int_a^b g(x)dx$$

برای این منظور توزیع دلخواهی بر $[a, b]$ با چگالی $f(x)$ اختیار می کنیم . علاوه بر متغیر تصادفی X ، که بر فاصله $[a, b]$ با چگالی $f(x)$ نعرف شده است احتیاج به متغیر تصادفی

$$h(x) = \frac{g(x)}{f(x)}$$

نیز نیاز دائم درین صورت

$$\theta = E(h(x)) = \int_a^b \frac{g(x)}{f(x)} f(x) dx$$

حال از چگالی f نمونه تصادفی X_1, X_2, \dots, X_n را استخراج می کنیم. آنگاه $h(X_1), h(X_2), \dots, h(X_n)$ نیز یک نمونه تصادفی بوده و داریم ،

$$E\left(\frac{1}{n} \sum_{i=1}^n h(X_i)\right) = \theta$$

طبق قانون قوی اعداد بزرگ

$$\frac{1}{n} \sum_{i=1}^n h(X_i) \xrightarrow{P} \theta$$

بنابراین $\frac{1}{n} \sum_{i=1}^n h(X_i)$ برآورد معقولی برای θ می باشد . واين برآورد معادل تقریب انتگرال مورد نظر است . با این روش بسیاری از انتگرال ها را که نمی توان با روش تحلیلی حل کرد تقریب زد .

مثال ۱: (برآورد π)

فرض کنید $(X, Y) \sim U(-1, 1)^2$ آنگاه از احتمال می دانیم که :

$$\theta = p(X^2 + Y^2 \leq 1) = \frac{\text{Area of } c}{\text{Area of } \Omega} = \frac{\pi}{4}$$

با تعریف

$$g(x, y) = I_C(x, y) = 1 \quad X^2 + Y^2 \leq 1$$

$$= \circ \quad o.w$$

داریم :

$$\theta = E(g(X, Y)) = \frac{\pi}{4}$$

برای تقریب θ از روش مونت کارلو به شیوه زیر عمل می کنیم .

گام ۱ : از توزیع $(1, 0)$ نمونه تصادفی $U^1, U^2, \dots, U_n^1, U_n^2, \dots, U_n^n$ را استخراج کرده و قرار می دهیم $1 - 2U_i^1 - 1 = 2U_i^2 - 1$ و $X_i = Y_i$ برای $i = 1, 2, \dots, n$ در این صورت $.Y_i \sim U(-1, 1)$ و $X_i \sim U(-1, 1)$

گام ۲ : $g(X_1, Y_1), \dots, g(X_n, Y_n)$ را محاسبه می کنیم .

گام ۳ : θ را با $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n g(X_i, Y_i)$ که نسبت تعداد (X_i, Y_i) هایی است که در دایره C قرار گرفته اند به n می باشد ، تقریب می زنیم . و در نهایت اگر بخواهیم π را تقریب بزنیم ، داریم

$$\pi \simeq 4\hat{\theta}$$

اگر علاوه ممند به محاسبه مساحت یا حجم غیر معمول D که معادله تحلیلی آن را در اختیار نداریم ، باشیم با تعریف $(\cdot) = I_D(\cdot)$ مشابه فوق به تقریب (برآورد) D می پردازیم .

به این منظور این ناحیه را در یک مستطیل ، مکعب یا حجره محصور کرده و از توزیع یکنواخت بر مستطیل ، مکعب یا حجره ، نمونه U_1, U_2, \dots, U_n را استخراج کرده و همچون مثال

قبل از $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n g(U_i)$ به عنوان تقریب D استفاده

می کنیم .

مثال ۲ :

فرض کنیم علاقمند به محاسبه $\frac{a(i)}{N}$ که در آن N بسیار بزرگ و $a(i)$ پیچیده بوده و محاسبه آن ساده نباشد ، باشیم . برای تقریب θ به روش مونت کارلو به قرار زیر عمل می کنیم .

گام ۱ : از توزیع $(1, 0)$ نمونه تصادفی $U^1, U^2, \dots, U_n^1, U_n^2, \dots, U_n^n$ را استخراج کرده ، $1 + X_i = [NU_i]$ را محاسبه می کنیم ، برای $i = 1, 2, \dots, n$. در این صورت X_i دارای توزیع یکنواخت بر $\{1, 2, \dots, N\}$ می باشد .

گام ۲ : $a(X_1), \dots, a(X_n)$ را محاسبه می کنیم .

گام ۳ : $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n a(i)$ را به عنوان برآورد θ در نظر می گیریم .

روشن است که روش مونت کارلو را برای توابع چند متغیره نیز می توان بکار برد . در واقع مزیت و برتری این روش به روش های عددی معمولی در حل مسائل با بعد بیشتر از یک نمایان می شود . زیرا روش های عددی معمولی را در حالت یک بعدی می توان به کار برد .

فرض کنیم $\mathfrak{R} \rightarrow \mathfrak{R}$: g و علاقمند به محاسبه انتگرال g بر ناحیه $A = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_d, b_d]$ باشیم .

$$\theta = \int_{a_1}^{b_1} \int_{a_2}^{b_2} \dots \int_{a_d}^{b_d} g(x_1, \dots, x_d) dx_1 \dots dx_d$$

برای این منظور چگالی دلخواه f بر A را در نظر می‌گیریم و تابع $h(\cdot)$ را به صورت

$$h(x_1, \dots, x_d) = \frac{g(x_1, \dots, x_d)}{f(x_1, \dots, x_d)}$$

تعريف می‌کنیم، در این صورت

$$\theta = E(h(X_1, \dots, X_d)) = \int_{a_1}^{b_1} \int_{a_2}^{b_2} \dots \int_{a_d}^{b_d} \frac{g(x_1, \dots, x_d)}{f(x_1, \dots, x_d)} f(x_1, \dots, x_d) dx_1 \dots dx_d$$

حال از چگالی (\cdot) نمونه تصادفی $\mathbf{x}_1, \dots, \mathbf{x}_d$ که در آن (X_{1i}, \dots, X_{di}) را استخراج می‌کنیم، آنگاه $h(\mathbf{X}_1, \dots, h(\mathbf{X}_n)$ یک نمونه تصادفی بوده و داریم که :

$$\frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i) \xrightarrow{P} \theta$$

بنابراین $\frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i)$ یک برآورد معقول برای θ یعنی انتگرال مورد نظر می‌باشد.

۳.۱ چگونگی انتخاب چگالی

در انتخاب f بدیهی است که باید f را انتخاب کرد که به سادگی قابل نمونه گیری باشد. زیرا اساس روش مونت کارلو استخراج نمونه و تقریب θ با استفاده از این نمونه می‌باشد. به هر حال برای هر f قابل نمونه گیری داریم

$$\theta = E\left(\frac{g(X)}{f(X)}\right) = E(h(X))$$

که در آن

$$h(X) = \frac{g(X)}{f(X)}$$

با این حال واریانس و درنتیجه خطای استاندارد برآورده θ بسته به انتخاب f می‌تواند کوچک یا بزرگ باشد که مطلوب آنست کمترین مقدار خطای استاندارد برآورده را داشته باشیم. یعنی

$$var(h(X)) = \int_a^b \left(\frac{g(x)}{f(x)} \right)^2 f(x) dx - \theta^2$$

کمینه شود. می‌توان نشان داد که وقتی $\frac{|g(x)|}{f(x)}$ ثابت باشد این مقدار کمینه می‌شود. و

$$f(x) = \frac{|g(x)|}{\int_a^b |g(x)| dx}$$

چگالی مطلوب می باشد .

مشاهده می شود که انتخاب بهترین چگالی در عمل به خاطر محاسبه انتگرال $\int_a^b |g(x)| dx$ که آسانتر از محاسبه $\int_a^b g(x)dx$ نمی باشد ، غیر ممکن است پس f را انتخاب می کنیم که به سادگی قابل نمونه گیری بوده و نسبت $\frac{|g(x)|}{f(x)}$ کراندار باشد. این امر باعث می شود که واریانس برآوردگر $\sum_{i=1}^n \frac{g(X_i)}{f(X_i)}$ کاهش یابد.

در بسیاری از مسائل چگالی f داده شده است ، به عنوان مثال در استنباط بیزی ، علاقمند به محاسبه انتگرال زیر هستیم

$$\gamma = E(h(\theta) | \mathbf{x}) = \int_{\Theta} h(\theta) \pi(\theta | \mathbf{x}) d\theta$$

و برای تقریب این انتگرال به روش مونت کارلو کافی است که نمونه‌ای از چگالی $(\mathbf{x} | \theta)$ استخراج کرده و از $\hat{\gamma} = \frac{1}{n} \sum_{i=1}^n h(\theta_i)$ برای تقریب γ استفاده می کنیم. بنابراین از این به بعد فرض می کنیم در پی محاسبه امید ریاضی تابعی از متغیر تصادفی X با توزیع f به روش مونت کارلو هستیم . یعنی

$$\gamma = E(h(X)) = \int_X h(x) f(x) dx$$

۴.۱ استخراج نمونه از چگالی f

حال اگر متغیر تصادفی X را با چگالی f در نظر بگیریم برای استخراج نمونه از این چگالی می توان از دو روش زیر استفاده کرد.

روش تبدیل معکوس تابع توزیع : می دانیم که برای متغیر تصادفی X داریم :

$$F(X) \sim U(0, 1)$$

حال اگر از توزیع $(1, 0) U$ نمونه تصادفی U_1, \dots, U_n را استخراج کنیم آنگاه با تبدیل

$$X_i = F^{-1}(U_i) = \inf \{x : F(x) \geq U_i\} \quad i = 1, \dots, n$$

X_1, \dots, X_n یک نمونه تصادفی از چگالی f می باشد .

مثال ۱:

فرض کنیم بخواهیم از چگالی $f(x) = \lambda e^{-\lambda x}$ نمونه‌ای تصادفی استخراج کنیم . بدین منظور :

گام ۱: U_1, \dots, U_n را از $U(0, 1)$ استخراج می‌کیم.

گام ۲: قرار می‌دهیم

$$X_i = F^{-1}(U_i) = \frac{-\log(1 - U_i)}{\lambda}$$

آنگاه X_1, \dots, X_n یک نمونه تصادفی از چگالی f می‌باشد.

مثال ۲:

برای استخراج نمونه از چگالی

$$f(x) = \frac{\lambda^m e^{-\lambda x} x^{m-1}}{\Gamma(m)}$$

که در آن $m \in \mathbb{N}$ و $\lambda > 0$ معلوم باشند. چون فرم بسته‌ای برایتابع توزیع این چگالی موجود نیست لذا از روش معکوس تابع توزیع نمی‌توان استفاده کرد. اما از روابط بین توزیعی توزیع گاما و توزیع نمایی، می‌توان نمونه ای به قرار زیر بدست آورد. اگر متغیرهای مستقل و هم توزیع با توزیع $E(\lambda)$ باشند آنگاه

$$X = \sum_{i=1}^m Y_i \sim \Gamma(m, \lambda)$$

گام ۱: نمونه تصادفی U_1, \dots, U_n را از $U(0, 1)$ استخراج می‌کیم.
گام ۲:

$$Y_i = \frac{-\log(1 - U_j^i)}{\lambda} \quad i = 1, \dots, m \quad j = 1, \dots, n$$

بدست می‌آوریم.

گام ۳: قرار می‌دهیم:

$$X_j^i = \sum_{i=1}^m Y_j^i$$

در این صورت $X_j \sim \Gamma(m, \lambda)$ ، بنابراین X_1, \dots, X_n یک نمونه تصادفی از چگالی

$$f(x) = \frac{\lambda^m e^{-\lambda x} x^{m-1}}{\Gamma(m)}$$

می‌باشد. مشابه این مثال موارد دیگر نیز پیش می‌آید که مستقیماً نمی‌توان از روش معکوس تابع توزیع استفاده کرد اما با استفاده از روابط بین توزیعی می‌توان این مشکل را

برطرف کرد. مثلا برای توزیع نرمال $N(\mu, \sigma^2)$ قادر به استفاده از روش تبدیل معکوس تابع توزیع نیستیم. اما می دانیم اگر U_1, U_2 متغیرهای تصادفی مستقل و هم توزیع از $(0, 1)$ باشند و قرار دهیم

$$X_1 = \mu + \sigma \sqrt{-2 \log U_1} \cos(2\pi U_2)$$

$$X_2 = \mu + \sigma \sqrt{-2 \log U_1} \sin(2\pi U_2)$$

آنگاه X_1, X_2 متغیرهای تصادفی مستقل و هم توزیع از توزیع $N(\mu, \sigma^2)$ خواهند بود.

۱.۴.۱ روش پذیرش – رد

فرض کنیم بخواهیم از چگالی f نمونه ای استخراج کنیم و این امر با روش تبدیل معکوس تابع توزیع و روابط بین توزیعی مقدور نباشد. اگر چگالی g با تکیه گاهی شامل تکیه گاه f چنان موجود باشد که

۱: به سادگی قابل نمونه گیری باشد .

۲: نسبت $\frac{f(x)}{g(x)}$ کراندار باشد یعنی یک $1 \leq M$ چنان وجود داشته باشد به طوری که $\frac{f(x)}{g(x)} \leq M \quad \forall x \in \chi$ آن گاه می توان با روش زیر نمونه مورد نظر را گرفت :

قرار می دهیم $e(x) = g(x)M \quad \forall x$. بنابراین $e(x) \geq f(x) \quad \forall x$ ، که تابع $e(x)$ را پوششی برای f می نامند. برای نمونه گیری به طریق زیر عمل می کنیم .

گام ۱: Y را از چگالی g استخراج می کنیم .

گام ۲: U را از $(0, 1)$ استخراج می کنیم .

گام ۳: اگر $U > \frac{f(Y)}{e(Y)}$ آنگاه Y را رد می کنیم . یعنی مقدار Y را به عنوان یک مشاهده از چگالی g در نظر نمی گیریم .

گام ۴: در غیر این صورت مقدار Y را پذیرفته و قرار میدهیم $Y = X$ ، یعنی X را به عنوان یک مشاهده از چگالی f در نظر می گیریم و برای بدست آوردن مشاهدات بعدی

گام های ۱ تا ۴ را تکرار می کنیم تا به نمونه مورد نظر دست یابیم .

لازم به ذکر است که نمونه بدست آمده با این روش دارای چگالی f می باشند. زیرا

$$P(X \leq y) = P(Y \leq y \mid U \leq \frac{f(Y)}{e(Y)})$$

$$= \frac{P(Y \leq y, U \leq \frac{f(Y)}{e(Y)})}{P(U \leq \frac{f(Y)}{e(Y)})}$$

$$\begin{aligned}
&= \frac{\int_{-\infty}^y \int_0^{\frac{f(z)}{e^{c(z)}}} g(z) dudz}{\int_{-\infty}^{\infty} \int_0^{\frac{f(z)}{e^{c(z)}}} g(z) dudz} \\
&= \int_{-\infty}^y f(z) dz
\end{aligned}$$

چند تذکر :

- ۱: شرط کرانداری $\frac{f}{g} \leq M$ شرط حیاتی است .
- ۲: احتمال پذیرش Y از استخراج شده از g ، $\frac{1}{M}$ می باشد . و این نشان می دهد که اگر M بزرگ باشد نرخ پذیرش $(\frac{1}{M})$ کوچک بوده والگوریتم کارا نخواهد بود .
- ۳: چگونه می توان M را کوچک ساخت ؟ بدیهی است که مقدار M به g انتخاب شده بستگی دارد . اگر g شبیه f باشد آنگاه احتیاج به M بزرگ نخواهد بود . پس بنابراین g ای انتخاب می کنیم که به سادگی قابل نمونه گیری بوده و هر چه بیشتر مشابه f باشد .
- ۴: اگر دم های چگالی g کوتاهتر از f باشد ، آنگاه M مناسبی نمی توان پیدا کرد .

مثال ۱ :

فرض کنیم می خواهیم از توزیع نرمال دم بریده در نقطه c ، نمونه ای به حجم n استخراج کنیم . با استفاده از روش پذیرش – رد

$$g(x) = \lambda e^{-\lambda x} \quad x > 0 \quad \lambda > 0$$

را پیشنهاد می کنیم . می دانیم که استخراج نمونه از چگالی نمایی با استفاده از روش تبدیل معکوس تابع توزیع براحتی امکان پذیر بوده و نیز دم توزیع نمایی ضخیم تراز توزیع نرمال است . بنابراین توزیع نمایی انتخاب مناسبی است . در اینجا

$$f(x) = \frac{\phi(x+c)}{1 - \Phi(c)}$$

که در آن

$$\Phi(c) = P(Z \leq c) \quad Z \sim N(0, 1)$$

$$\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

حال می خواهیم $1 \geq M$ ای را پیدا کنیم که $\frac{f}{g}$ به آن کراندار بوده یعنی

$$Mg(x) \geq f(x) \quad \forall x > c$$

و همچنین M کمترین مقدار را اختیار کند . بدین منظور قرار می دهیم

$$M = \sup_{\{x: x > c\}} \frac{f(x)}{g(x)}$$

$$= \sup_{x > c} \frac{\phi(x + c)}{\lambda e^{-\lambda x} (1 - \Phi(c))}$$

$$= \frac{1}{\lambda \sqrt{2\pi} (1 - \Phi(c))} \sup_{x > c} \exp\left\{ \lambda x - \frac{(x + c)^2}{2} \right\}$$

سوپریمم در نقطه $x = \lambda - c$ رخ می دهد . و بنابراین :

$$M = \frac{e^{\left(\frac{\lambda^2 - 2\lambda c}{2}\right)}}{\sqrt{2\pi} \lambda (1 - \Phi(c))}$$

حال می توان λ را طوری انتخاب کرد که M کمینه و درنتیجه $\frac{1}{M}$ بیشینه گردد .
می توان نشان داد که

$$\lambda = \frac{c + \sqrt{c^2 + 4}}{2}$$

مقدار مطلوب برای λ است .

۲.۴.۱ محدودیت های روش پذیرش رد

دیدیم که روش پذیرش - رد ابراز قدرتمندی برای استخراج نمونه از توزیع های نچندان خوش رفتار می باشد . اما مواردی نیز پیش می آید که قادر به استفاده از این روش نیستیم .
از این موارد به دو مورد زیر می توان اشاره کرد .

۱ : شایعترین مورد زمانی اتفاق می افتد که چگالی هدف ، $f(x)$ به طور کامل معلوم نباشد و تنها می دانیم که $f(x) = c\tilde{f}(x)$. که در آن c مجھول می باشد . در استنباط بیزی به کرات با این مورد برخورد می کنیم .

۲ : پیدا کردنتابع چگالی w که به سادگی قابل نمونه گیری باشد و نسبت $\frac{f}{w}$ کراندار بماند، همیشه کار چندان ساده ای هم نیست . به علاوه ممکن است که چنین w ای را پیدا کنیم اما کران M عدد نسبتا بزرگی باشد و درنتیجه احتمال پذیرش $\frac{1}{M}$ ، ناچیز بوده و این کارایی روش پذیرش - رد را به شدت کاهش می دهد

فصل ۲

زنجیره‌های مارکف

۱.۲ تعاریف

دنباله‌ی متغیرهای تصادفی $\{X^{(t)} : t \geq \circ\}$ ، که همه‌ی $X^{(t)}$ ها در مجموعه‌ی شمارای S مقدار می‌گیرند، را «زنجیره‌ی مارکف همگن» با فضای وضعیت S گوییم هرگاه در شرط زیر صدق کند:

$$\begin{aligned} Pr(X^{(t+1)} = j | X^{(t)} = i, X^{(t-1)} = i_{t-1}, \dots, X^{(\circ)} = i_{\circ}, X^{(\circ)} = i_{\circ}) \\ = Pr(X^{(t+1)} = j | X^{(t)} = i) = p_{ij} \end{aligned}$$

p_{ij} را احتمال انتقال یک مرحله‌ای از وضعیت i به وضعیت j گفته، و قرار می‌دهیم: $P = [p_{ij}]_{i,j \in S}$ در اینصورت P ماتریس احتمال انتقالات یک مرحله‌ای زنجیره‌ی $\{X^{(t)}\}$ است. بدیهی است که: $\forall i, j \in S: p_{ij} \geq \circ; \sum_{j \in S} p_{ij} = ۱$ یعنی هر سطر ماتریس P یک احتمال شرطی بر فضای وضعیت S تعریف می‌کند.

مشابهًا اگر تعریف کنیم: $P^{(m)} = [p_{ij}^{(m)}]_{i,j \in S}; p_{ij}^{(m)} = Pr(X^{t+m} = j | X^t = i)$ آنگاه $p_{ij}^{(m)}$ را احتمال انتقال m مرحله‌ای از وضعیت i به وضعیت j ، و $P^{(m)}$ را ماتریس احتمال انتقال m مرحله‌ای زنجیره‌ی $\{X^{(t)}\}$ می‌گوییم. از آنجایی که:

$$\begin{aligned} p_{ij}^{(\circ)} &= Pr(X^{(\circ)} = j | X^{(\circ)} = i) = \sum_{k \in S} Pr(X^{(\circ)} = j, X^{(\circ)} = k | X^{(\circ)} = i) \\ &= \sum_{k \in S} Pr(X^{(\circ)} = j | X^{(\circ)} = k) Pr(X^{(\circ)} = k | X^{(\circ)} = i) = \sum_{k \in S} p_{ik} p_{kj} \end{aligned}$$

. $P^{(\circ)} = P \times P$:

به همین ترتیب می توان نشان داد :

$$P^{(m)} = \underbrace{P \times \dots \times P}_m$$

$$P^{(m)} \times P^{(l)} = P^{(m+l)}$$

که به برابری های « چپمن - کلموگروف » مشهورند و بیان می کنند که با داشتن ماتریس احتمال انتقالات یک مرحله‌ای، ماتریس احتمال انتقالات چند مرحله‌ای را می توان بدست آورد. لذا ماتریس احتمال انتقالات یک مرحله‌ای (P) کاملاً ساختار انتقالات زنجیره‌ی مارکف $\{X^{(t)}\}$ را مشخص می کند. به عبارت دیگر: هر زنجیره مارکف با ماتریس احتمال انتقالات یک مرحله‌ای اش قابل بیان است (با داشتن P زنجیره $\{X^{(t)}\}$ را داریم).

۲.۲ ساختار احتمالی زنجیره

$X^{(\circ)}$ را وضعیت اولیه زنجیره می گویند و با داشتن توزیع تمام $X^{(\circ)}$ ها را خواهیم داشت. زیرا :

فرض کنیم $X^{(\circ)}$ دارای توزیع $\Pi^{(\circ)}$ باشد. این توزیع را - برای سادگی نمادگذاری - بصورت زیر نمایش می دهیم :

$$\pi_i^{(\circ)} = Pr(X^{(\circ)} = i)$$

$$\Pi^{(\circ)} = [\pi_i^{(\circ)}]_{i \in S}$$

حال اگر توزیع $X^{(1)}$ را با $\Pi^{(1)}$ نشان دهیم در این صورت :

$$\Pi^{(1)} = [\pi_j^{(1)}]_{j \in S}$$

$$\pi_j^{(1)} = Pr(X^{(1)} = j) = \sum_{i \in S} Pr(X^{(1)} = j \mid X^{(\circ)} = i) Pr(X^{(\circ)} = i)$$

$$= \sum_{i \in S} \pi_i P_{ij}$$

با نماد گذاری ماتریسی یعنی :

$$(\Pi^{(1)})^T = (\Pi^{(\circ)})^T \times P$$

مشابهًاً اگر توزیع $X^{(t)}$ را با $\Pi^{(t)}$ نشان دهیم داریم :

$$\Pi^{(t)} = [\pi_j^{(t)}]_{j \in S}$$

$$(\Pi^{(t)})^T = (\Pi^{(t)})^T \times P$$

که نشان می دهد توزیع $X^{(t)}$ تنها به ماتریس P و توزیع $\Pi^{(\circ)}$ بستگی دارد. به عبارت دیگر P و $\Pi^{(\circ)}$ کل ساختار احتمالی زنجیره را بدست می دهند.

۳.۲ توزیع حدی زنجیره

دیدیم که برای زنجیره مارکف با ماتریس احتمال انتقالات یک مرحله‌ای P ، با انتخاب یک توزیع احتمال برای وضعیت اولیه زنجیره $(X^{(0)})$ ، دنباله‌ی توزیع های احتمال $\{\Pi^{(t)}\}_{t \geq 0}$ برای $X^{(t)}$ ها بدست آمد. برای این دنباله از توزیع ها چنانچه داشته باشیم :

$$\lim_{t \rightarrow \infty} \Pi^{(t)} = \Pi$$

و Π یک توزیع احتمال بر S باشد، آنگاه Π را توزیع حدی زنجیره گوییم.
روشن است که برای انتخاب های مختلف $(\Pi^{(0)})$ ، توزیع های حدی در صورت وجود - ممکن است متفاوت باشند. آیا زنجیره‌هایی هستند که برای هر انتخاب $(\Pi^{(0)})$ ، توزیع حدی همواره وجود داشته و یکتا باشد ؟
در تعریف زیر این نوع زنجیره‌ها را نامگذاری کرده و به ویژگی مطلوب آنها برای استخراج نمونه اشاره می کنیم. در بخش های بعد با ارائه مقدمات و سپس قضیه ای (قضیه ارگودیک)، شرط‌های کافی برای این نوع زنجیره‌های مارکف را می آوریم.
زنجیره‌ی $\{X^{(t)}\}$ را «ارگودیک» گوییم هر گاه برای هر توزیع ابتدایی $(\Pi^{(0)})$ ، دنباله توزیع های احتمال بدست آمده $\{\Pi^{(t)}\}_{t \geq 0}$ به توزیع احتمال Π همگرا باشد. یعنی :

$$\forall \Pi^{(0)} : \lim_{t \rightarrow \infty} \Pi^{(t)} = \Pi$$

به عبارت دیگر برای زنجیره ارگودیک، صرف نظر از توزیع اولیه زنجیره $(\Pi^{(0)})$ ، توزیع $X^{(t)}$ ها $(\Pi^{(t)})$ به توزیع حدی یکتای Π میل می کند.
همین مسئله یک ایسیده برای استخراج نمونه به مامی دهد : فرض کنیم بخواهیم از چگالی $f(x)$ با تکیه گاه شمارای S نمونه‌ای استخراج کنیم اما با روش‌های معمول ^۱ قادر به این کار نباشیم. اگر زنجیره‌ی مارکف ارگودیکی داشته باشیم که توزیع حدی یکتای آن $f(x)$ باشد و از هر سطر ماتریس احتمال انتقالات یک مرحله‌ای آن (P) بتوان نمونه تولید کرد ^۲ به روش زیر عمل می کنیم:

(۱) از توزیع دلخواه $(\Pi^{(0)})$ ، $(X^{(0)})$ را تولید می کنیم.

(۲) با داشتن $X^{(t)} = x^{(t)}$ - بدست آمده از مرحله‌ی قبل - از سطر $x^{(t)}$ ام ماتریس P ، $X^{(t+1)}$ را تولید می کنیم.

(۳) این کار را تا بدست آمدن تعداد مشاهدات مورد نیاز ادامه می دهیم.

در اینصورت با توجه به ارگودیک بودن زنجیره، توزیع $X^{(t)}$ ها برای t های به قدر کافی بزرگ به $f(x)$ میل می کند. پس می توان نتیجه گرفت که نمونه بدست آمده از الگوریتم

^۱ روش معکوس تابع توزیع و روش پذیرش - رد

^۲ هر سطر ماتریس احتمال انتقالات، یک توزیع شرطی بر فضای وضعیت S تعریف می کند.

فوق یک نمونه تقریبی از $f(x)$ است. برای بهبود کیفیت این نمونه (یعنی نزدیکی بیشتر توزیع $X^{(t)}$ به $f(x)$) و برای کاهش تأثیر انتخاب توزیع اولیه نمونه بسیار بزرگی به روش فوق استخراج کرده، تنها مشاهدات انتهایی بدست آمده را در نظر بگیریم و از مشاهدات ابتدایی صرف نظر کنیم.

نتایج فوق نتها بر پایه ای ارگودیک بودن زنجیره بناشده بود. حال باید بینیم شرایط ارگودیک بودن یک زنجیره مارکف چیست؟

۴.۲ توزیع ایستای زنجیره

اگر برای وضعیت اولیه زنجیره مارکف، $(X^{(0)}, [\lambda_i]_{i \in S})$ ، توزیع $\Lambda = [\lambda_i]_{i \in S}$ را که برای آن داریم:

$$(\Lambda)^T = (\Lambda)^T \times P$$

در نظر بگیریم، آنگاه برای هر $t \geq 0$ توزیع $X^{(t)}$ نیز $\Lambda = [\lambda_i]_{i \in S}$ خواهد بود. در این حالت Λ را توزیع ایستای زنجیره $\{X^{(t)}\}$ (یا متناظرًا توزیع ایستای ماتریس P) می‌گوییم. واضح است که با انتخاب Λ به عنوان توزیع اولیه زنجیره Λ توزیع حدی زنجیره نیز می‌باشد.

۵.۲ زنجیره‌های تحويلی ناپذیر و نامتناوب

برای هر $i, j \in S$ ، گوییم وضعیت i را قابل دسترسی است اگر:

$$\exists k \in N : p_{ij}^{(k)} = P(X^{(t+k)} = j | X^{(t)} = i) > 0$$

به عبارت دیگر از وضعیت i پس از مدت زمان متناهی بتوان به وضعیت j رفت. دو وضعیت i و j را مرتبط گوییم هرگاه زاز i قابل دسترسی و i از j قابل دسترسی باشد. زنجیره‌ی مارکف $\{X^{(t)}\}$ را تحولی ناپذیر گوییم هرگاه تمام وضعیت‌های آن مرتبط باشند. برای $i \in S$ دوره‌ی تناوب وضعیت i را با $d(i)$ نشان داده و آنرا به صورت زیر تعریف می‌کنیم:

$$d(i) = \gcd\{k \in N : p_{ii}^{(k)} > 0\}$$

که در آن \gcd نشان دهنده‌ی بزرگترین مقسوم علیه مشترک یک مجموعه از اعداد می‌باشد. اگر $1 = d(i)$ وضعیت i را نامتناوب می‌گوییم. به عبارت دیگر وضعیت i را

نامتناوب است هرگاه از π در هر زمانی بتوان دوباره به π برگشت. بطور کلی زنجیره را نا متناوب گوییم هرگاه هر وضعیت آن نامتناوب باشد.

قضیه ارگودیک: برای زنجیره مارکف تحویل ناپذیر و نامتناوب $\{X^{(t)}\}$ ، اگر بتوان یک جواب غیر صفر برای معادله

$$(Z)^T \times P = (Z)^T ; \quad Z = [z_i]_{i \in S}$$

بدست آورد طوری که: $\sum_{i \in S} z_i < \infty$ ، آنگاه زنجیره $\{X^{(t)}\}$ ارگودیک است.
نتیجه: زنجیره مارکف تحویل ناپذیر و نامتناوب $\{X^{(t)}\}$ اگر توزیع ایستا داشته باشد آنگاه ارگودیک است. در این حالت این توزیع ایستا، توزیع حدی یکتای زنجیره مارکف $\{X^{(t)}\}$ می باشد.

پس اگر یک زنجیره مارکف تحویل ناپذیر و نامتناوب باشد، برای ارگودیک بودن آن کافی است نشان دهیم که توزیع ایستا دارد و با تعیین آن، توزیع ایستای حدی یکتای آن بدست می آید. اما نشان دادن این مطلب مستلزم پیدا کردن توزیع $\Pi = [\pi_i]_{i \in S}$ است که برای آن داشته باشیم: $(\Pi)^T \times P = (\Pi)^T$ که جز در موارد خیلی خاص کار مشکل، و عمدهاً غیرممکن است. ویرگی زمان برگشت پذیری کار را ساده‌تر می کند.

۶.۲ زنجیره زمان برگشت پذیر

زنجیره مارکف $\{X^{(t)}\}$ را زمان برگشت پذیر گویند هرگاه یک توزیع بر S یافت شود که برای آن داشته باشیم:

$$\forall i, j \in S : \pi_i p_{ij} = \pi_j p_{ji}$$

اگر زنجیره مارکف $\{X^{(t)}\}$ به معنی فوق زمان برگشت پذیر باشد آنگاه Π توزیع ایستای زنجیره بوده ویکتاست. زیرا، با جمع بر روی i دوطرف تساوی فوق داریم:

$$\forall j \in S : \sum_{i \in S} \pi_i p_{ij} = \pi_j$$

و این یعنی: $(\Pi)^T = (\Pi)^T \times P$.

نتیجه: زنجیره مارکف تحویل ناپذیر، نامتناوب و زمان برگشت پذیر ارگودیک است.

۷.۲ زنجیره مارکف با فضای وضعیت ناشمارا

دیدیم که با ساختن یک زنجیره مارکف ارگودیک که توزیع ایستای حدی یکتاش چگالی هدف $f(x)$ باشد، می توان یک نمونه تقریبی از $f(x)$ استخراج کرد. بدین منظور در این

بخش فضای وضعیت زنجیره را از حالت شمارا به حالت ناشمارا تعمیم داده، و تعاریف و نتایج بدست آمده در بخش‌های قبل را در این حالت بررسی می‌کنیم.
دنباله‌ی متغیرهای تصادفی $\{X^{(t)} : t \geq 0\}$ ، که همه‌ی $X^{(t)}$ ها^۳ در مجموعه‌ی ناشمارای S مقدار می‌گیرند، را «**زنジیره مارکف همگن**» با فضای وضعیت S گوییم هرگاه در شرط زیر صدق کند:

$$\begin{aligned} \forall A \in \mathcal{F}, \forall t \geq 0, \forall x_0, x_1, \dots, x_{t-1}, x \in S &: \\ \Pr(X^{(t+1)} \in A | X^{(t)} = x, X^{(t-1)} = x_{t-1}, \dots, X^{(1)} = x_1, X^{(0)} = x_0) \\ &= \Pr(X^{(t+1)} \in A | X^{(t)} = x) = \int_A k(x, y) dy \end{aligned}$$

که در آن \mathcal{F} یک σ -میدان بر S ، و $k(x, y)$ برای $x \in S$ ثابت، چگالی توزیع شرطی K بر S است (یعنی $(Y | X = x) \sim K$). به عبارت دیگر $\{k(x, \cdot) : x \in S\}$ خانسواره‌ای از چگالی‌های شرطی بر S می‌باشد. در این حالت $k(x, y)$ را هسته‌ی احتمال انتقالات^۴ یک مرحله‌ای زنجیره گوییم و همان نقش ماتریس احتمال انتقالات یک مرحله‌ای را در اینجا ایفا می‌کند.
اگر برای وضعیت اولیه‌ی $(x^{(0)}, X^{(0)})$ ، توزیع $\pi^{(0)}$ با چگالی $\pi^{(0)}(x)$ را در نظر بگیریم، آنگاه مانند زنجیره مارکف با فضای وضعیت شمارا، $\pi^{(0)}$ کل ساختار احتمالی زنجیره را مشخص می‌کنند. زیرا برای هر $t \geq 0$ داریم :

اگر $X^{(t)}$ داری توزیع $\Pi^{(t)}$ با چگالی $\pi^{(t)}(x)$ باشد و $X^{(t+1)}$ داری توزیع $\Pi^{(t+1)}$ با چگالی $\pi^{(t+1)}(x)$ ، آنگاه :

$$\begin{aligned} P(X^{(t+1)} \in A) &= \int_A \pi^{(t+1)}(y) dy = \int_A \int_S \pi^{(t)}(x) k(x, y) dx dy \\ \pi^{(t+1)}(y) &= \int_S \pi^{(t)}(x) k(x, y) dx \end{aligned}$$

و این یعنی توزیع $X^{(t+1)}$ با $\pi^{(t+1)}$ قابل بیان است. پس برای زنجیره مارکف با هسته‌ی احتمال انتقالات یک مرحله‌ای $k(x, y)$ ، با انتخاب یک توزیع احتمال $\Pi^{(0)}$ برای وضعیت اولیه زنجیره $(x^{(0)}, X^{(0)})$ ، دنباله‌ی توزیع های احتمال $\{\Pi^{(t)}\}_{t \geq 0}$ با دنباله‌ی چگالی‌های متناظر $\{\pi^{(t)}(x)\}_{t \geq 0}$ برای $X^{(t)}$ ها بدست می‌آید. در اینصورت اگر دنباله‌ی توابع چگالی $\{\pi^{(t)}(x)\}_{t \geq 0}$ بر S به تابع π همگرای یگنواخت باشد، که در آن π یک تابع چگالی بر S است، آنگاه $(x^{(0)}, \pi)$ را چگالی حدی زنجیره، و توزیع Π متناظر با چگالی π را توزیع حدی زنجیره گوییم. در اینجا نیز برای انتخاب های مختلف $\Pi^{(0)}$ توزیع های حدی-در صورت وجود—ممکن است متفاوت باشند و تنها برای زنجیره‌های ارگودیک، صرف نظر از توزیع اولیه زنجیره، توزیع $X^{(t)}$ ها به توزیع حدی یکنای Π میل می‌کند.

^۳ فرض بر این است که همه‌ی $X^{(t)}$ ها پیوسته‌ی مطلق اند، یعنی توزیع آنها چگالی دارد.
^۴ Transition kernel: با ماتریس احتمال انتقالات مقایسه کنید.

اگر برای توزیع Λ با چگالی $\lambda(x)$ داشته باشیم :

$$\forall y \in S : \lambda(y) = \int_S \lambda(x) k(x, y) dx$$

در اینصورت Λ را توزیع ایستا، و $\lambda(x)$ را چگالی ایستای زنجیره می‌گوییم. واضح است که با انتخاب Λ به عنوان توزیع اولیه زنجیره، Λ توزیع حدی زنجیره نیز می‌باشد.

تمام تعاریف و نتایج ارگودیک بودن، تحويل ناپذیری، نامتناوب بودن و زمان بازگشتی بودن مشابه با آنچه در حالت زنجیره با فضای وضعیت شمارا آمده است می‌باشد، تفاوت تنها در استفاده از $k(x, y)$ به جای p_{ij} ، $\pi^{(t)}(x)$ به جای $\pi_i^{(t)}$ ، و انتگرال به جای سیگما می‌باشد.

فصل ۳

روش مونت کارلو با زنجیره‌های مارکف ($MCMC$)

دیدیم که برای تقریب کمیت

$$\theta = E(X) = \int_S h(x)f(x)dx$$

به روش مونت کارلسو، با استخراج مشاهدات مستقل و هم توزیع، از چگالی $f(x)$ و استفاده از $\hat{\theta}_{MC} = \frac{1}{n} \sum_{i=1}^n h(X_i)$ کمیت θ را تقریب می‌زیم. اما استخراج نمونه از چگالی $f(x)$ مشکل اساسی روش مونت کارلو می‌باشد. برای این کار دو روش معکوس تابع توزیع و روش پذیرش و رد را معرفی کرده، و محدودیتهای آنها را بر شمردیم. اکنون به دنبال روش قدرتمندتری برای استخراج نمونه هستیم که محدودیتهای کمتری داشته باشد.

در بحث زنجیره‌های مارکف ارگودیک، الگوریتمی را معرفی کردیم که به کمک آن می‌توان از چگالی هدف $f(x)$ با تکیه گاه S نمونه استخراج کرد. برای این کار، لازم بود که یک زنجیره‌ی مارکف ارگودیک با توزیع ایستای حدی یکتای $f(x)$ داشته باشیم بطوریکه از هر سطر ماتریس احتمال انتقالات در حالت S شمارا و یا از هر چگالی (y, k) ، برای $x \in S$ ثابت، در حالت S ناشمارا بتوان نمونه تولید کرد. به علاوه دیدیم که زنجیره‌ی مارکف تحويل ناپذیر، نامتناوب و زمان برگشت پذیر ارگودیک است. حال مسئله این است که چگونه می‌توان یک زنجیره‌ی مارکف تحويل ناپذیر، نامتناوب و زمان برگشت پذیر با توزیع ایستای حدی یکتای $f(x)$ بر S ساخت طوریکه از سطرهای ماتریس احتمال انتقالات و یا چگالیهای هسته احتمال انتقالات بتوان نمونه تولید کرد.

در واقع روش مونت کارلو با زنجیره‌های مارکف ($MCMC$)، چنین زنجیره‌ای را برای استخراج یک نمونه‌ی تقریبی از $f(x)$ ساخته (بخش زنجیره‌ی مارکف روش) و با استخراج

نمونه، روش مونت کارلو (بخش مونت کارلو روش) را برای تقریب کمیتها را به کار برد. مارکف نشان داد که قانون اعداد بزرگ در حالت وابستگی مارکفی نمونه، باز هم برقرار است. بنابراین در بخش مونت کارلو روش، چنانچه از تقریبی بودن توزیع کناری مشاهدات بدست آمده چشم پوشی، وفرض کنیم که مشاهدات ما وابسته‌ی مارکفی‌اند و توزیع کناری آنها $f(x)$ است آنگاه با پشتونه‌ی قانون اعداد بزرگ می‌توان روش مونت کارلو را برای تقریب کمیت‌ها به کار برد. لذا بخش مونت کارلو روش با پذیرش یک تقریب سر راست است. اما مسئله‌ی اصلی ساختن زنجیره‌ی مارکف تحويل ناپذیر، نامتناوب و زمان برگشت پذیر با چگالی ایستای حدی یکتای (x^f) ، یعنی بخش زنجیره‌ی مارکف $MCMC$ می‌باشد.

برای تولید چنین زنجیره‌ای دو الگوریتم ارائه شده است : الگوریتم متروپلیس-هستینگ^۱ و نمونه‌گیری گیبس^۲.

۱۰.۳ الگوریتم متروپلیس-هستینگ

این الگوریتم با انتخاب یک هسته‌ی احتمال انتقالات $(x | y)$ بر S ، با داشتن مشاهده $X^{(t)}$ ، پس از سه گام، مشاهدی $X^{(t+1)}$ را به ما می‌دهد. بدین منظور به یک مقدار اولیه برای بدست آوردن سایر مشاهدات نیاز داریم :

در زمان $t = 0$ $x^{(0)} = x^*$ را از توزیع ابتدایی دلخواه $\Pi^{(0)}$ با چگالی $(x^{(0)}, \pi^{(0)})$ ، که بر S تعریف شده است، با این شرط که $x^{(0)} > f(x^{(0)})$ ، تولید کند. البته می‌توان $x^{(0)} \in S$ با شرط $x^{(0)} > f(x^{(0)})$ ، به دلخواه انتخاب کرد که معادل است با انتخاب یک توزیع بر S با چگالی $\pi^{(0)}(x^{(0)}) = 1$ ؛ $\forall y \in S, y \neq x^{(0)}$: $\pi^{(0)}(y) = 0$.

حال با داشتن $x^{(t)}$ ، الگوریتم مقدار $X^{(t+1)} = x^{(t+1)}$ را بصورت زیر تولید می‌کند:

(۱) مقدار پیشنهادی X^* را از چگالی $(x^{(t)}, g(x^{(t)})$ تولید می‌کیم.

(۲) نسبت $R(u, v) = \frac{f(v)g(u|v)}{f(u)g(v|u)}$ را که بصورت $R(u, v)$ تعریف شده است را برای $(x^{(t)}, X^*)$ محاسبه می‌کنیم.

(۳) $X^{(t+1)}$ را بصورت زیر انتخاب می‌کنیم :

• با احتمال $\min\{1, R(x^{(t)}, X^*)\}$

• با احتمال $1 - \min\{1, R(x^{(t)}, X^*)\}$

Metropolis-Hasting	۱
Gibbs Sampling	۲

این الگوریتم را تا زمانی که تعداد مشاهدات مورد نیاز را بدست آمد، تکرار می کیم. لذا این الگوریتم با انتخاب هسته‌ی احتمال انتقالات $(x | y)$ بر S ، و مقدار اولیه‌ی $x^{(0)}$ یک زنجیره ما می دهد.

توجه

- مقدار X^* را که در گام یک تولید می شود، مقدار پیشنهادی می گوییم چون در گام سوم پذیرفته یا رد می شود.
- چگالی شرطی $(y | x)$ را که با آن مقادیر پیشنهادی X^* را با داشتن $x^{(t)}$ تولید می کنیم، چگالی پیشنهادی می گوییم.
- نسبت متراپلیس-هستینگ $R(x^{(t)}, X^*)$ همواره تعریف می شود. زیرا مقدار پیشنهادی $X^* = x^{(t)}$ تنها وقتی تولید می شود که $f(x^{(t)}) > g(x^{(*)} | x^{(t)}) > 0$.
- مارکف بودن زنجیره بدست آمده با این روش بدیهی است زیرا $X^{(t+1)}$ تنها به $X^{(t)}$ بستگی دارد.
- تحويل ناپذیری و نامتناوب بودن زنجیره به انتخاب چگالی پیشنهادی $(y | x)$ بستگی دارد. بنابراین باید $(y | x)$ طوری انتخاب شود که درنهایت زنجیره بودست آمده از الگوریتم تحويل ناپذیر و نامتناوب باشد.
- زمان برگشت پذیری زنجیره همواره برقرار است.

زیرا با انتخاب $f(x)$ به عنوان توزیع اولیه برای تولید $X^{(*)}$ داریم : اگر $(X^{(*)}, X^{(1)})$ باشد آنگاه :

$$w(x^{(*)}, x^{(1)}) = Pr(X^{(1)} = x^{(1)} | X^{(*)} = x^{(*)})f(x^{(*)})$$

ولی، برای احتمال شرطی داریم :

$$\begin{aligned} & Pr(X^{(1)} = x^{(1)} | X^{(*)} = x^{(*)}) \\ &= Pr(X^* = x^{(1)}, U \leq \min\{1, R(x^{(*)}, X^*)\} | X^{(*)} = x^{(*)}) \\ &= Pr(U \leq \min\{1, R(x^{(*)}, X^*)\} | X^* = x^{(1)}, X^{(*)} = x^{(*)}) \\ &\quad Pr(X^* = x^{(1)} | X^{(*)} = x^{(*)}) \\ &= Pr(U \leq \min\{1, R(x^{(*)}, x^{(1)})\} \times g(x^{(1)} | x^{(*)})) \end{aligned}$$

اگر $R(x^{(*)}, x^{(1)}) \leq 1$: آنگاه $f(x^{(1)})g(x^{(*)} | x^{(1)}) \leq f(x^{(*)})g(x^{(1)} | x^{(*)})$

$$\begin{aligned} & Pr(X^{(1)} = x^{(1)} | X^{(*)} = x^{(*)}) = R(x^{(*)}, x^{(1)})g(x^{(1)} | x^{(*)}) \\ &= \frac{f(x^{(1)})g(x^{(*)} | x^{(1)})}{f(x^{(*)})g(x^{(1)} | x^{(*)})}g(x^{(1)} | x^{(*)}) = \frac{f(x^{(1)})g(x^{(*)} | x^{(1)})}{f(x^{(*)})} \end{aligned}$$

و اگر $R(x^{(0)}, x^{(1)}) > 1$ آنگاه $f(x^{(1)})g(x^{(0)} | x^{(1)}) > f(x^{(0)})g(x^{(1)} | x^{(0)})$

$$Pr(X^{(1)} = x^{(1)} | X^{(0)} = x^{(0)}) = g(x^{(1)} | x^{(0)})$$

پس :

$$\begin{aligned} w(x^{(0)}, x^{(1)}) &= f(x^{(1)})g(x^{(0)} | x^{(1)}) ; R(x^{(0)}, x^{(1)}) \leq 1 \\ &= f(x^{(0)})g(x^{(1)} | x^{(0)}) ; R(x^{(0)}, x^{(1)}) > 1 \end{aligned}$$

مشابه‌ها می‌توان نشان داد:

$$\begin{aligned} w(x^{(1)}, x^{(0)}) &= f(x^{(1)})g(x^{(0)} | x^{(1)}) ; R(x^{(0)}, x^{(1)}) \leq 1 \\ &= f(x^{(0)})g(x^{(1)} | x^{(0)}) ; R(x^{(0)}, x^{(1)}) > 1 \end{aligned}$$

يعنى $w(x^{(0)}, x^{(1)}) = w(x^{(1)}, x^{(0)})$ ، به عبارت دیگر چکالی توان $(X^{(0)}, X^{(1)})$ متقارن است. بنابراین توزيع کساري $X^{(0)}$ و $X^{(1)}$ یکسان و همان $f(x)$ می‌باشد. این مطلب نه تنها برای $(X^{(0)}, X^{(1)})$ بلکه برای هر $(X^{(t)}, X^{(t+1)})$ صادق است. لذا :

$$\begin{aligned} \forall x, y \in S : Pr(X^{(t)} = x)Pr(X^{(t+1)} = y | X^{(t)} = x) &= \\ Pr(X^{(t+1)} = y)Pr(X^{(t)} = x | X^{(t+1)} = y) \end{aligned}$$

يعنى :

$$\forall x, y \in S : f(x)g(y | x) = f(y)g(x | y)$$

و اين يعني زنجيره‌ی بدست آمده زمان برگشت پذير است، بنابراین $f(x)$ توزيع ايستای يكتای زنجيره می‌باشد.^۳

در اگوريتم متروپليس-هستينگ مسئله‌ی اصلی پيدا کردن $(x | y)$ -ایي است که اولاً برای هر $x \in S$ ثابت قابل نمونه گيري باشد و ثانياً باعث شود زنجيره‌ی حاصل از الگوريتم ارگوديك باشد. علاوه بر اينها، ملاكه‌ای دیگري را نيز برای انتخاب $(x | y)$ مناسب باید مد نظر قرار داد. واضح است که مشاهدات حاصل از الگوريتم باید نمونه‌ی خوبی از جامعه‌ی $f(x)$ باشند. به عنوان مثال زنجيره‌ای که بيشتر مقادير پيشنهادی را رد می‌کند، منجر به مشاهداتی از $f(x)$ می‌شود که خيلي از آنها تكراري و ثابت هستند و بنابراین اين نمونه نمی‌تواند پراكندگی جامعه‌ی $f(x)$ را به خوبی نشان دهد. همچنين زنجيره‌ای که همه‌ی مقادير پيشنهادی را پذيرد، منجر به مشاهداتی از $f(x)$ می‌شود که بسيار پراكنده بوده و اين در سرعت همگرايی توزيع مشاهدات به توزيع حدی $f(x)$ تأثير

^۳ اثبات برای حالت پيوسته مشابه همين است.

گذاشته، باعث می شود برای بدست آورد یک نمونه‌ی قابل قبول، مشاهدات بسیار زیادی را با الگوریتم تولید کنیم. لذا انتخاب $(y | x)$ مناسب در الگوریتم متروپلیس-هستینگ حیاتی است. اجراهای عملی این الگوریتم نشان داده که با انتخاب $(x | y)$ مناسب، نتایج بسیار رضایت‌بخش‌اند.

اما متأسفانه در حالت کلی نمی توان قاعده‌ای برای انتخاب $(x | y)$ مناسب ارائه کرد. بلکه باید با توجه به ساختار چگالی $f(x)$ و تکیه گاه آن S ، و با درنظر گرفتن شرط‌ها و ملاک‌های انتخاب $(x | y)$ مناسب، $(y | x)$ را انتخاب کرد. بهترین کار دسته‌بندی مسائلی است که با روش $MCMC$ در پی حل آنها هستیم. در آمار بیشترین کاربرد این روش در استنباط بیزی است. وقتی که ما می خواهیم میانگین‌پسین، میانه‌ی پسین، احتمال‌های پسین و ... را محاسبه کنیم، اما در حل معمول انتگرال (بعض‌سیگما) حاصل از تعریف این کمیتها با مشکل رویروییم. استفاده از زنجیره‌های مستقل حداقل در این مورد خاص^۴ بسیار کارساز است.

۱.۱.۳ الگوریتم متروپلیس-هستینگ و زنجیره‌های مستقل

اگر در الگوریتم متروپلیس-هستینگ، $(y | x)$ طوری انتخاب شود که داشته باشیم : $(y | x) = g(y)$ ، که در آن $(x | y)$ چگالی یک توزیع بر S است، در اینصورت هر مقدار پیشنهادی X^* مستقل از مشاهده‌ی قبلی $x^{(t)}$ از $(x | y)$ تولید می شود. لذا زنجیره‌ی حاصل از الگوریتم مستقل می باشد.

با این انتخاب، نسبت متروپلیس-هستینگ به فرم زیر است :

$$R(x^{(t)}, X^*) = \frac{f(X^*)g(x^{(t)})}{f(x^{(t)})g(X^*)}$$

به علاوه اگر برای هر $x \in S$ که $f(x) > 0$ داشته باشیم : $0 < g(x) < \infty$ آنگاه زنجیره حتماً تحويل ناپذیر و نامتناوب بوده، لذا ارجو دیک با چگالی ایستای حدی یکتای $f(x)$ می باشد. بنابراین در این حالت دیگر نیازی به چک کردن تحويل ناپذیری و نامتناوب بودن زنجیره‌ی حاصل از الگوریتم نبوده، بلکه کافیست شرط فوق در مورد $(x | y)$ صدق کند.

مثال : استنباط بیزی

در استنباط بیزی با مسئله‌ی زیر مواجه‌ایم :

مشاهدات y_1, \dots, y_k را با تابع درستنمایی $(y | \theta)$ داریم. به علاوه، برای پارامتر θ چگالی پیشین $(\theta | \pi)$ را در نظر گرفته ایم. استنباط بیزی (برآورد نقطه‌ای یا فاصله‌ای و آزمون فرض) در مورد θ تماماً برپایه‌ی تابع چگالی پسین زیر می باشد :

^۴ البته این مورد آنقدرها هم خاص نیست!

$$\begin{aligned}\forall \theta \in \Theta & : \pi(\theta | y) \propto \pi(\theta)L(\theta | y) \\ \forall \theta \in \Theta & : \pi(\theta | y) = c\pi(\theta)L(\theta | y)\end{aligned}$$

که ضریب تناسب بصورت زیر محاسبه می شود :

$$c = (\int_{\Theta} \pi(\theta)L(\theta | y)d\theta)^{-1}$$

جز در موارد استاندارد، معمولاً محاسبه ای انتگرالی که c را تعیین می کند پیچیده است.
بنابراین در استنباط بیز معمولاً c را نداریم.

در مسئله‌ی برآورد بیزی θ ، چنانچه تابع زیان مسئله را مربع خطأ در نظر بگیریم،
برآورد زیر را برای پارامتر مجھول θ داریم:

$$\hat{\theta}_B = \int_{\Theta} \theta \pi(\theta | y) d\theta = \int_{\Theta} \theta c\pi(\theta)L(\theta | y)d\theta$$

از آنجا که محاسبه‌ی c را دشوار فرض کردیم، دلیلی وجود ندارد که محاسبه‌ی $\hat{\theta}_B$ را ساده‌تر از محاسبه‌ی c فرض کنیم. بنابراین به نظر می رسد برای داشتن برآورد بیزی باید به تقریب $\hat{\theta}_B$ قانع شد. نیز فرض کنیم که $\pi(\theta)$ قابل نمونه گیری بوده و می خواهیم با روش مونت کارلو این کار را انجام دهیم.

برای این کار ممکن است این راه حل به ذهن مان برسد : ابتدا با استخراج نمونه از $\pi(\theta)$ ، روش مونت کارلو را برای تقریب c به کار ببریم. سپس با استخراج نمونه‌ای دیگر از $\pi(\theta)$ ، $\hat{\theta}_B$ را با $\frac{1}{n} \sum_{i=1}^n c\theta_i L(\theta_i | y)$ تقریب بزیم. به این راه حل دو ایجاد وارد است : اولاً ما برای تقریب $\hat{\theta}_B$ از دوبار تقریب، تقریب c و بعد تقریب $\hat{\theta}_B$ با روش مونت کارلو، استفاده کرده ایم که مسلماً دقت تقریب $\hat{\theta}_B$ را کاهش میدهد. دوماً معمولاً محاسبه‌ی $L(\theta | y)$ به دلیل فرم حاصل ضربی آن ساده نبوده و اغلب با تقریب و خطاهای گرد کردن همراه است. بنابراین این تقریب نیز بر تقریب‌های اشاره شده قبلی اضافه شده، دقت برآورد را کاهش می دهد.

راه حل منطقی استفاده از روش مونت کارلو با زنجیره‌های مارکف (*MCMC*) می باشد : با انتخاب $\pi(\theta^{(t)} | \theta^{(t)}) = \pi(\theta^* | \theta^{(t)})$ واستفاده از الگوریتم متروپلیس-هستینگ، یک زنجیره مستقل با توزیع ایستای حدی یکتای $\pi(\theta | y)$ بدست می آوریم (زیرا برای هر $\theta \in \Theta$ که $\pi(\theta | y) > 0$ ، به دلیل متناسب بودن، داریم $\pi(\theta) > 0$). آنگاه با داشتن مشاهدات کافی (y_1, \dots, y_m) از زنجیره و حذف مشاهدات اولیه m باید بهبود کیفیت نمونه و کاهش تاثیر مقدار اولیه، $\hat{\theta}_B$ را با $\frac{1}{n-m} \sum_{t=m+1}^n \theta^{(t)}$ تقریب می زیم.

مثال : برآورد پارامتر آمیخته
فرض کنیم y_1, \dots, y_{100} را بطور مستقل و همتوزیع از توزیع آمیخته‌ی

$$\delta N(7, \frac{1}{3}) + (1 - \delta)N(10, \frac{1}{3}) \quad ; \quad \delta \in (0, 1)$$

مشاهده کرده‌ایم. به عنوان مثال این مشاهدات ممکن است مربوط به دو جامعه‌ی $N(7, \frac{1}{4})$ و $N(10, \frac{1}{4})$ باشند که با احتمال مجهول δ از $N(7, \frac{1}{4})$ و احتمال $\delta - 1$ از $N(10, \frac{1}{4})$ می‌باشند. با فرض چگالی پیشین $(\pi(\delta) = 1 ; \delta \in [0, 1])$ ، هدف برآورد بیز پارامتر مجهول δ می‌باشد.تابع درستنمای مشاهدات عبارتست از:

$$L(\delta | y) = \prod_{j=1}^{100} \frac{1}{\sqrt{2\pi}} (\delta \exp(-2(y_j - 7)^2) + (1 - \delta) \exp(-2(y_j - 10)^2))$$

لذا چگالی پسین δ بصورت $cL(\delta | y) = cL(\delta | y)$ می‌باشد که در آن :

$$\begin{aligned} c^{-1} &= \int_0^1 L(\delta | y) d\delta \\ &= \int_0^1 \prod_{j=1}^{100} \frac{1}{\sqrt{2\pi}} (\delta \exp(-2(y_j - 7)^2) + (1 - \delta) \exp(-2(y_j - 10)^2)) \end{aligned}$$

با اینکه بعد δ یک است و مجموعه‌ای که انتگرال بر آن گرفته می‌شود خوشرفتار است ($[0, 1]$ در \mathbb{R} فشرده است) با این حال محاسبه‌ی c کار چندان راحتی به نظر نمی‌رسد! و اگر برآورد بیز δ را بخواهیم با انتگرال زیر مواجه‌ایم :

$$\int_0^1 \delta \pi(\delta | y) d\delta$$

اما $(\pi(\delta | y))$ به سادگی قابل نمونه گیری بوده و با توجه به آنچه در مثال قبل گفته شد استفاده از روش مونت کارلو با زنجیره‌های مارکف سر راست است. با انتخاب مقدار اولیه‌ی $\delta^{(0)}$ از توزیع یکنواخت بر $(0, 1)$ ، و انتخاب هسته‌ی احتمال انتقالات $(x | y) = \pi(x | y)$ بر $(0, 1)$ ، برای این مثال، الگوریتم متropolیس-هستینگ با داشتن $R(\delta^{(t)}, \delta^{(t+1)})$ را بصورت زیر تولید می‌کند :

(۱) δ^* را از یکنواخت بر $(0, 1)$ تولید می‌کنیم.

(۲) نسبت متropolیس-هستینگ $R(\delta^{(t)}, \delta^*) = \frac{\pi(\delta^* | y)}{\pi(\delta^{(t)} | y)}$ را محاسبه می‌کنیم.

(۳) $\delta^{(t+1)}$ را بصورت زیر انتخاب می‌کنیم :

$$\delta^{(t+1)} = \delta^* ; \min\{1, R(\delta^{(t)}, \delta^*)\}$$

$$= \delta^{(t)} ; 1 - \min\{1, R(\delta^{(t)}, \delta^*)\}$$

همانطور که ملاحظه شد، استفاده از زنجیره‌های مستقل تا حدی ساده بوده، به ویژه در استنباط بیزی بسیار راه گشاست. البته در استنباط بیزی اگر از توزیع پیشین نتوان به سادگی نمونه گرفت این روش کارساز نیست و باید به دنبال زنجیره‌های دیگری بود.

معمولًاً زنجیره‌های غیر مستقل سرعت همگرایی بیشتری نسبت به زنجیره‌های مستقل داشته، درنتیجه نمونه‌هایی با کیفیت تری تولید می‌کنند. بنابراین در مسائلی که می‌توان از آنها استفاده، بر زنجیره‌های مستقل ارجحیت دارند. علاوه بر زنجیره‌های مستقل زنجیره‌های ساده‌ی دیگری (مانند زنجیره‌ی گام تصادفی و ...) نیز برای انتخاب توزیع پیشنهادی دسته بندی شده‌اند. برای اطلاع از زنجیره‌های دیگر، [۱] منابعی را ذکر کرده است. حال دو مثال را می‌آوریم که در آنها از زنجیره‌های مستقل استفاده نشده است.

مثال : از جانسون و آلبرت ۱۹۹۹

کارخانه‌ای یکی از محصولاتش را بهبود داده، می‌خواهد نسبت مشتریانی که محصول جدید را به محصول قدیم ترجیح می‌دهند برآورد کند. این کارخانه مطمئن است که این نسبت از $\frac{1}{2}$ بیشتر است، یعنی $[1, \frac{1}{2}] \in p$. با انتخاب توزیع یکنواخت بر $[1, \frac{1}{2}]$ به عنوان توزیع پیشین پارامتر p ، برآورد بیز p مد نظر است. یک بررسی نمونه‌ای نشان داده، از ۲۰ مشتری متقاضی محصول، ۱۲ نفر محصول جدید را ترجیح داده‌اند. چگالی پیشین،تابع درستنمایی و چگالی پسین برای این مسئله بصورت زیر است :

$$\begin{aligned}\pi(p) &= 2 \quad ; \quad p \in [\frac{1}{2}, 1] \\ L(p | y) &= p^{12}(1-p)^8 \\ \pi(p | y) &= cp^{12}(1-p)^8 \quad ; \quad p \in [\frac{1}{2}, 1], \quad c^{-1} = \int_{\frac{1}{2}}^1 p^{12}(1-p)^8 dp\end{aligned}$$

برآورد بیز p از حل $\int_0^1 cp^{13}(1-p)^7 dp$ بدست می‌آید. علاقمندیم که این انتگرال را به روش مونت کارلو با زنجیره‌های مارکف، منتهایه با زنجیره‌ی مستقل حل کنیم. باید $(p^* | p^{(t)})$ مناسی را که در شرطها و ملکهای گفته شده صدق کند پیدا کنیم. در مورد $p \in [\frac{1}{2}, 1]$ این کار بسیار سخت است. راه حل جانسون و آلبرت بدین صورت است : با دوره پارامتر کردن مسئله هم فضای پارامتر را می‌توان طوری تعییر داد که انتخابهای بیشتری برای توزیع پیشنهادی داشته باشیم و هم فرم توزیع پسین انعطاف پذیرتر می‌شود.^۵ بدین منظور تبدیل $\theta = \ln(\frac{p - \frac{1}{2}}{1 - p})$; $\theta \in \mathbb{R}$ که عکس این تبدیل $p = \frac{\frac{1}{2} + \exp(\theta)}{1 + \exp(\theta)}$ است، را برای پارامتر p در نظر می‌گیریم. در اینصورت چگالی پسین به فرم زیر است :

$$\pi(\theta | y) \propto \frac{(\frac{1}{2} + \exp(\theta))^{12} \exp(\theta)}{(1 + \exp(\theta))^{22}} \quad ; \quad \theta \in \mathbb{R}$$

حال برای θ انتخاب $(\theta^*, \theta^{(t)})$ مناسب ساده‌تر است. جانسون و آلبرت چگالی پیشنهادی زیر را در نظر گرفته‌اند :

^۵ علاوه بر این موارد، دوباره پارامتره کردن مناسب مدل مزایای دیگری مانند : کاهش واریانس برآورده، افزایش سرعت اجرای الگوریتم و ...، نیز دارد. به مرجع [۱] رجوع شود.

$$g(\theta^*, \theta^{(t)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(\theta^* - \theta^{(t)})^2)$$

که چگالی توزیع $N(\theta^{(t)}, \sigma^2)$ است و σ^2 پیش از اجرای الگوریتم تعیین می شود. σ^2 در مناسب بودن توزیع پیشنهادی، کیفیت مشاهدات زنجیره و سرعت همگرایی به توزیع حدی نقش مهمی دارد. با این انتخاب و انتخاب مقدار اولیه زنجیره $(\theta^{(0)}, N(0, 1))$ از توزیع الگوریتم متropolیس-هستینگ با داشتن $(\theta^{(t+1)}, \theta^{(t)})$ را بصورت زیر تولید می کند:

(1) θ^* را از توزیع $N(\theta^{(t)}, \sigma^2)$ تولید می کنیم.

(2) نسبت متropolیس-هستینگ $R(\theta^{(t)}, \theta^*) = \frac{\pi(\theta^*|y)g(\theta^{(t)}|\theta^*)}{\pi(\theta^{(t)}|y)g(\theta^*|\theta^{(t)})}$ را محاسبه می کنیم.

(3) $\theta^{(t+1)}$ را بصورت زیر انتخاب می کیم :

$$\begin{aligned} \theta^{(t+1)} &= \theta^* ; \quad \min\{1, R(\theta^{(t)}, \theta^*)\} \\ &= \theta^* ; \quad 1 - \min\{1, R(\theta^{(t)}, \theta^*)\} \end{aligned}$$

مثال : برآورد بیز پارامترهای توزیع وایبل توزیع وایبل با پارامترهای α و λ ، دارای چگالی به فرم زیر است :

$$f(x | \alpha, \lambda) = \alpha \lambda x^{\alpha-1} \exp(-\lambda x^\alpha) ; \quad x > 0, \alpha, \lambda > 0$$

این توزیع عضو خانواده توزیعهای نمایی نمی باشد و با انتخاب یک توزیع پیشین برای پارامترهای (α, λ) ، توزیعهای پسین دقیقی برای (α, λ) بدست نمی دهد. لذا برآورد بیز (α, λ) با مشکل مواجه است. سعی می کنیم با روش مونت کارلو با زنجیرهای مارکف، برآورد بیزی برای (α, λ) پیدا کنیم.

فرض کنیم مشاهدات $(x_1, x_2, x_3) = (\frac{2}{10}, \frac{1}{10}, \frac{25}{100})$ را از توزیع وایبل مشاهده کرده ایم و توزیع پیشین $\pi(\alpha, \lambda) \propto \exp(-\alpha) \lambda^{\beta-1} \exp(-\beta \lambda)$ را برای (α, λ) در نظر گرفته ایم که β ثابت و معلوم است.تابع درستنمایی و چگالی پسین برای این مسئله بصورت زیر است :

$$\begin{aligned} L(\alpha, \lambda | x) &= \left(\frac{5}{100}\right)^{\alpha-1} \alpha^3 \lambda^3 \exp(-\lambda((\frac{1}{10})^\alpha + (\frac{2}{10})^\alpha + (\frac{25}{100})^\alpha)) \\ \pi(\alpha, \lambda | x) &\propto \left(\frac{5}{100}\right)^{\alpha-1} \alpha^3 \exp(-\alpha) \lambda^{\beta-2} \exp(-\lambda((\frac{1}{10})^\alpha + (\frac{2}{10})^\alpha + (\frac{25}{100})^\alpha + \beta)) \\ c^{-1} &= \end{aligned}$$

$$\int_0^\infty \int_0^\infty \left(\frac{5}{100}\right)^{\alpha-1} \alpha^3 \exp(-\alpha) \lambda^{\beta-2} \exp(-\lambda((\frac{1}{10})^\alpha + (\frac{2}{10})^\alpha + (\frac{25}{100})^\alpha + \beta)) d\alpha d\lambda$$

حل دو انتگرال زیر برآورد بیز (α, λ) را به ما می دهد :

$$\begin{aligned} \int_0^\infty \int_0^\infty \left(\frac{5}{100}\right)^{\alpha-1} \alpha^4 \exp(-\alpha) \lambda^{\beta-2} \exp(-\lambda((\frac{1}{10})^\alpha + (\frac{2}{10})^\alpha + (\frac{25}{100})^\alpha + \beta)) d\alpha d\lambda \\ \int_0^\infty \int_0^\infty \left(\frac{5}{100}\right)^{\alpha-1} \alpha^3 \exp(-\alpha) \lambda^{\beta-1} \exp(-\lambda((\frac{1}{10})^\alpha + (\frac{2}{10})^\alpha + (\frac{25}{100})^\alpha + \beta)) d\alpha d\lambda \end{aligned}$$

روشن است که این انتگرالها را نمی توان به روش‌های معمول حل کرد لذا حل به روش مونت کارلو با زنجیره‌های مارکف گزینه‌ی مناسبی است. با انتخاب چگالی پیشنهادی (α^*, λ^*) که حاصل‌ضرب چگالی دو توزیع نمایی با میانگین‌های $\alpha^{(t)}$ و $\lambda^{(t)}$ است، و انتخاب مقدارهای اولیه‌ی $\alpha^{(0)}$ و $\lambda^{(0)}$ از توزیع نمایی با میانگین ۱، الگوریتم متropolیس-هستینگ با داشتن $(\alpha(t), \lambda^{(t)})$ ، $(\alpha^{(t+1)}, \lambda^{(t+1)})$ را بصورت زیر تولید می کند :

(۱) α^* را از توزیع نمایی با میانگین $(\alpha^{(t)}, \lambda^{(t)})$ و λ^* را از توزیع نمایی با میانگین $\lambda^{(t)}$ تولید کرده، تا (α^*, λ^*) را بدست آوریم.

۲ نسبت متropolیس-هستینگ

$$R((\alpha^{(t)}, \lambda^{(t)}), (\alpha^*, \lambda^*)) = \frac{\pi((\alpha^*, \lambda^*)|x)g((\alpha^{(t)}, \lambda^{(t)})|(\alpha^*, \lambda^*))}{\pi((\alpha^{(t)}, \lambda^{(t)})|x)g((\alpha^*, \lambda^*)|(\alpha^{(t)}, \lambda^{(t)}))}$$

را محاسبه می کنیم.

(۳) $(\alpha^{(t+1)}, \lambda^{(t+1)})$ را بصورت زیر انتخاب می کنیم :

$$(\alpha^{(t+1)}, \lambda^{(t+1)}) = (\alpha^*, \lambda^*) \quad ; \quad \min\{1, R((\alpha^{(t)}, \lambda^{(t)}), (\alpha^*, \lambda^*))\}$$

$$= (\alpha^{(t)}, \lambda^{(t)}) \quad ; \quad 1 - \min\{1, R((\alpha^{(t)}, \lambda^{(t)}), (\alpha^*, \lambda^*))\}$$

دو نمودار که معمولاً برای بررسی زنجیره‌های تولید شده با الگوریتم متropolیس-هستینگ مورد استفاده قرار می گیرد، سیر نمونه و هیستوگرام مشاهدات زنجیره می باشند. در نمودار مسیر نمونه زنجیره، مشاهدات $X^{(t)}$ در مقابل زمان t دریک محور مختصات رسم می شوند. با این نمودار سرعت و وضعیت همگرایی زنجیره، تأثیر مقدار اولیه و نحوه‌ی گردش و تغییرات مشاهدات زنجیره در S را می توان دید و در مورد کارایی زنجیره بدنست آمده، و در نتیجه الگوریتم قضاوت کرد. همچنین هیستوگرام مشاهدات زنجیره نیز می تواند ملاکی برای کارایی زنجیره تولید شده با الگوریتم باشد. مسلماً زنجیره‌ای کارتر است که هیستوگرام مشاهدات آن به نمودار چگالی هدف $f(x)$ شبیه‌تر باشد.

مانند سایر الگوریتم‌های عددی، چند بار اجرا کردن الگوریتم و بررسی نتایج این اجرایها نیز ملاک دیگری برای کارایی الگوریتم می باشد. اگر در چند اجرای مختلف، الگوریتم نتایج تقریباً یکسانی به ما داد، نتیجه می گیریم الگوریتم کاراست. در غیر این صورت، نمی توان به نتایج الگوریتم زیاد اعتماد کرد.

۲.۳ نمونه گیری گیبز

برای تقریب $\theta = E(X) = \int_S h(x)f(x)dx$ به روش مونت کارلو وقتی که با روش‌های معمول نتوان نمونه‌ای از $f(x)$ استخراج کرد، الگوریتم متروپلیس–هستینگ روشی عمومی برای تولید یک زنجیره ارگودیک با چگالی ایستای حدی یکتاوی $f(x)$ می‌باشد که تحقق‌های این زنجیره را می‌توان به عنوان نمونه‌ای تقریبی از $f(x)$ فرض کرد. اما اجرای این الگوریتم زمانی امکان پذیر است که هسته‌ی احتمال انتقالات (چگالی پیشنهادی) مناسبی را برای تولید یک مشاهده به شرط داشتن مشاهده قبلی، داشته باشیم. اشاره شد که بطور کلی نمی‌توان قاعده‌ای برای انتخاب توزیع پیشنهادی ارائه کرد.

بنابراین بهتر است براساس مسائلی که به روش $MCMC$ می‌خواهیم حل کنیم، الگوریتم‌های مخصوصی را طرح کرده و زنجیره‌های تولید شده با این الگوریتمها را دسته بندی کنیم. مثلاً در الگوریتم متروپلیس–هستینگ، یک دسته بندی براساس انواع زنجیره‌هایی که می‌توان با این الگوریتم تولید کرد، زنجیره‌های مستقل است.

حال با طرح مسئله‌ای، الگوریتم دیگری را برای حل این مسئله طرح می‌کنیم:
فرض کنیم $(X_1, \dots, X_p) \sim f(\mathbf{x})$. $\mathbf{X}_{-i} = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_p)$ ؛ $i = 1, \dots, p$

(بردار (X_1, \dots, X_p) با حذف مولفه‌ی i آن). همچنین برای $i = 1, \dots, p$ ، توزیع شرطی $X_i | \mathbf{X}_{-i} = \mathbf{x}_{-i}$ با چگالی $f(X_i | \mathbf{x}_{-i})$ را درنظر می‌گیریم.
می‌خواهیم $\theta = E(\mathbf{X}) = \int_S h(\mathbf{x})f(\mathbf{x})d\mathbf{x}$ را به روش مونت کارلو تقریب بزنیم، اما به روش‌های معمول نمی‌توانیم از (\mathbf{x}) نمونه استخراج کنیم. چنانچه برای $i = 1, \dots, p$ ، $f(x_i | \mathbf{x}_{-i})$ طوری باشد که با داشتن \mathbf{x}_{-i} بتوان از آن نمونه تولید کرد، آنگاه با انتخاب وضعیت اولیه‌ی (\mathbf{x}^0) ، الگوریتم زیر با داشتن $\mathbf{X}^{(t+1)}$ را تولید می‌کند:

(۱) یک ترتیب از مولفه‌های $\mathbf{x}^{(t)}$ را انتخاب می‌کنیم.

(۲) برای $i = 1, \dots, p$ ، X_i^* را از $f(x_i | \mathbf{x}_{-i}^{(t)})$ تولید کرده، تا بدست آید.

(۳) قرار می‌دهیم: $\mathbf{X}^{(t+1)} = \mathbf{X}^*$

توجه

- این الگوریتم را گیبز برای تقریب یک انتگرال در فیزیک پیشنهاد کرده است و به همین خاطر به نمونه گیری گیبز یا الگوریتم گیبز مشهور است.

- زنجبهای که با این الگوریتم بدست می آید، $\{X^{(t)}\}_{t \geq 0}$ ، مارکف است زیرا $X^{(t+1)}$ تنها به $X^{(t)}$ وابسته است.

در گام دوم به جای اینکه X^* از $f(x)$ تولید شود، با تولید p مولفه‌ی X_i^* بردار X^* بدست می آید. برای اینکه ترتیب تولید مولفه‌ها بر تولید X^* تأثیر نداشته باشد، در گام اول از بین p جایگشت ممکن از مولفه‌های $X^{(t)}$ یکی را انتخاب و گام دوم را براساس این ترتیب انتخاب شده اجرا می کنیم. به همین دلیل در برخی از موارد ابتدا یک توزیع را احتمال بر p جایگشت مولفه‌های $X^{(t)}$ در نظر گرفته، گام اول را چنین بیان می کنند: یک جایگشت را از توزیع تعریف شده بر جایگشتهای مولفه‌های $X^{(t)}$ تولید می کنیم. سپس گام دوم را با این ترتیب تولید شده اجرا می کنند.

- اجرای گام دوم بر اساس ترتیب انتخاب شده در گام اول را یک سیکل می گوییم.
- برای یک سیکل خاص، در i امین مرحله‌ی سیکل الگوریتم با داشتن $X^{(t)}$ ، مقدار پیشنهادی $(g_i(x^* | x^{(t)}) = (x_1^{(t)}, \dots, X_i^*, \dots, x_p^{(t)})$ را از چگالی پیشنهادی $(X^* = (x_1^*, \dots, X_i^*, \dots, x_p^*))$ تولید می کند. که در آن :

$$g_i(x^* | x^{(t)}) = f(x_i^* | x_{-i}^{(t)}) ; \quad x_{-i}^* = x_{-i}^{(t)}$$

سایر نقاط

در این حالت، $R(x^{(t)}, x^*) = \frac{f(x^*) g_i(x^{(t)} | x^*)}{f(x^{(t)}) g_i(x^* | x^{(t)})} \equiv 1$ یعنی مقدار پیشنهادی X^* پذیرفته می شود. بنابراین هر سیکل از p الگوریتم متropolیس-هستینگ تشکیل شده است و می توان الگوریتم فوق را حالت خاصی از متropolیس-هستینگ انگاشت.

- می دانیم که $f(x_i | x_{-i}) = c_i f(x)$ ، که در آن $c_i = \int f(x) dx_{-i}$. چنانچه در این الگوریتم $f(x_i | x_{-i})$ برای برخی از i ها نرمال نشده باشد \Rightarrow یعنی c_i مجھول باشد، چون هر سیکل در الگوریتم متشكل از p الگوریتم متropolیس-هستینگ است، مشکلی برای اجرای الگوریتم نخواهیم داشت.

- برای بهبود عملکرد این الگوریتم، در یک سیکل وقتی که مولفه‌ی اول پیشنهادی X_1^* با داشتن $X^{(t)}$ تولید شد، مولفه‌ی دوم پیشنهادی X_2^* را از $X^{(t+1)}$ به جای $(g_i(x^* | x_1^{(t)}, x_2^{(t)}, \dots, x_p^{(t)})$ تولید کرده و به همین ترتیب برای سایر مولفه‌های تا کامل شدن سیکل از مقدار جدید $x_i^{(t+1)}$ به جای $x_i^{(t)}$ استفاده می کنیم.

^۶ تابع چگالی است اگر $\int_{\mathbb{R}} f(x) dx = 1$ (۱) $\forall x : f(x) \geq ۰$ (۲). اگر f در شرط اول صدق کند و بدایم $\int_{\mathbb{R}} f(x) dx < \infty$ آنگاه f را یک چگالی نرمال نشده می گوییم.

- زنجیره‌ی مارکفی که با این الگوریتم تولید می‌شود تحويل ناپذیر، نامتناوب و زمان برگشت پذیر و درنتیجه ارگودیک بوده، بنابراین این الگوریتم نیز یکی از بزارهای روش $MCMC$ است.

مثال : استنباط بیزی

فرض کنیم $\theta = (\theta_1, \dots, \theta_p)$ پارامترهای یک مدل آماری باشند. همانطور که در بخش الگوریتم متropolیس-هستینگ اشاره شد، برای استنباط بیزی در مورد θ با توزیع پیشین $\pi(\theta)$ ، باید کمیتهای پسین را محاسبه کرد و یک روش مناسب استفاده از روش مونت کارلو با زنجیره‌های مارکف است. دیدیم که اگر $\pi(\theta)$ قابل نمونه‌گیری باشد، با انتخاب $\pi(\theta)$ به عنوان توزیع پیشنهادی یک زنجیره‌ی مستقل برای روش مونت کارلو با زنجیره‌های مارکف به دست می‌آید. اما اگر $\pi(\theta)$ قابل نمونه‌گیری نباشد چه؟

فرض کنیم $\pi(\theta)$ قابل نمونه‌گیری نبوده اما برای $p, i = 1, \dots, p$ $\pi(\theta_i | \theta_{-i})$ قابل نمونه‌گیری باشد، از الگوریتم گیز برای داشتن یک نمونه‌ی تقریبی از $(y | \theta)$ استفاده کرد.

مثال : مدل نرمال هایپر پارامتر^۷

فرض کنیم y_1, \dots, y_m را از $N(\mu, \frac{1}{\tau})$ مشاهده کرده‌ایم و معتقدیم که $\mu \sim G(2, 1)$ و $\tau \sim N(0, 1)$ در اینصورت :

$$f(\mathbf{y} | \mu, \tau) = (2\pi)^{-\frac{m}{2}} \tau^{\frac{m}{2}} \exp(-\frac{\tau}{2} \sum_{j=1}^m (y_j - \mu)^2)$$

$$\pi(\tau) = \tau \exp(-\tau) \quad \pi(\mu) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{\mu^2}{2})$$

بنابراین :

$$f(\mathbf{y}, \mu, \tau) = (2\pi)^{-\frac{m+1}{2}} \tau^{\frac{m+1}{2}} \exp(-\frac{\tau}{2} \sum_{j=1}^m (y_j - \mu)^2 - \frac{\mu^2}{2} - \tau)$$

داریم : $\pi(\mu, \tau | \mathbf{y}) \propto f(\mathbf{y}, \mu, \tau)$. برای داشتن نمونه از $(\mu, \tau | \mathbf{y})$ با روش‌های معمول نمی‌توان عمل کرد. اما $\pi(\mu | \tau, \mathbf{y})$ قابل نمونه‌گیری است :

$$\pi(\mu | \tau, \mathbf{y}) = \frac{f(\mathbf{y}, \mu, \tau)}{\pi(\tau, \mathbf{y})} = c f(\mathbf{y}, \mu, \tau) ; \quad c^{-1} = \int \int f(\mathbf{y}, \mu, \tau) d\tau d\mu$$

⁷ : Hyperparameter Model

در مدل‌های بیز معمول، فرض می‌شود که پارامترهای اصلی مدل متغیر تصادفی‌اند و دارای توزیع احتمال با پارامترهای ثابت می‌باشند. این پارامترها را اصطلاحاً پارامترهای فرعی مدل می‌گویند. برای بیان تجربیات و اعقایدات خود در قالب یک توزیع برای پارامترهای اصلی مدل، تعیین پارامترهای فرعی نقش بسیار مهمی دارد. چنانچه این پارامترهای فرعی را نیز متغیر تصادفی با توزیع احتمال مشخصی فرض کنیم، مدل را هایپر پارامتر می‌گوییم. مدل‌های هایپر پارامتر انعطاف پیشتری داشته ولی محاسبات عددی در آنها با پیچیدگی‌های زیادی روبرو است. استفاده از روش مونت کارلو با زنجیره‌های مارکف به ویژه الگوریتم گیز برای انجام محاسبات این مدلها، بسیار رایج است.

یعنی :

$$\pi(\mu | \tau, \mathbf{y}) \propto \exp\left(-\frac{1+m\tau}{\tau}(\mu - \frac{\tau \sum_{j=1}^m y_j}{1+m\tau})^2\right)$$

و یا $\mu | \tau, \mathbf{y} \sim N\left(\frac{\tau \sum_{j=1}^m y_j}{1+m\tau}, \frac{1}{1+m\tau}\right)$ داریم

$$\pi(\tau | \mu, \mathbf{y}) \propto \tau^{\frac{m}{2}+1} \exp\left(-\tau(1 + \frac{1}{\tau} \sum_{j=1}^m (y_j - \mu)^2)\right)$$

و یا $\tau | \mu, \mathbf{y} \sim G\left(\frac{m}{2} + 2, 1 + \frac{1}{\tau} \sum_{j=1}^m (y_j - \mu)^2\right)$

با انتخاب یک مقدار اولیه‌ی $(\mu^{(0)}, \tau^{(0)})$ ، آگوریتم گیبز با داشتن $(\mu^{(t)}, \tau^{(t)})$ در سه گام $(\mu^{(t+1)}, \tau^{(t+1)})$ را بصورت زیر تولید می‌کند :

۱) یکی از دو ترتیب (μ, τ) یا (τ, μ) را انتخاب می‌کنیم.

۲) با توجه به ترتیب انتخاب شده در گام اول، μ^* را از $N\left(\frac{\tau^{(t)} \sum_{j=1}^m y_j}{1+m\tau^{(t)}}, \frac{1}{1+m\tau^{(t)}}\right)$ و τ^* را از $G\left(\frac{m}{2} + 2, 1 + \frac{1}{\tau^{(t)}} \sum_{j=1}^m (y_j - \mu^{(t)})^2\right)$ تولید می‌کنیم.

$$(\mu^{(t+1)}, \tau^{(t+1)}) = (\mu^*, \tau^*) \quad (3)$$

مثال : مدل مولفه‌های واریانس

فرض کنیم برای $i = 1, \dots, k$ و $j = 1, \dots, m$ Y_{ij} ها مستقل و بوده و $Y_{ij} \sim N(\theta_i, \sigma_e^2)$. برای پارامترهای اصلی مدل، توزیع‌های پیشین زیر را در نظر می‌گیریم :

$a_1, b_1 \sim IG(a_1, b_1)$ و $\theta_1, \dots, \theta_k \sim N(\mu, \sigma_\theta^2)$ ثابت بوده اما $a_2, b_2, \sigma_e^2 \in \mathbb{R}^+$ و $\mu_0 \in \mathbb{R}$. $\sigma_\theta^2 \sim IG(a_2, b_2)$ و $\mu \sim N(\mu_0, \sigma_\mu^2)$ همگی معلوم‌اند.

$$f(\mathbf{y} | \theta_1, \dots, \theta_k, \mu, \sigma_\theta^2, \sigma_e^2) = (2\pi\sigma_e^2)^{-\frac{km}{2}} \exp\left(-\frac{1}{2\sigma_e^2} \sum_{i=1}^k \sum_{j=1}^m (y_{ij} - \theta_i)^2\right)$$

$$\pi(\theta_1, \dots, \theta_k | \mu, \sigma_\theta^2, \sigma_e^2) = (2\pi\sigma_\theta^2)^{-\frac{k}{2}} \exp\left(-\frac{1}{2\sigma_\theta^2} \sum_{i=1}^k (\theta_i - \mu)^2\right)$$

$$\pi(\mu | \theta_1, \dots, \theta_k, \sigma_\theta^2, \sigma_e^2) = (2\pi\sigma_\mu^2)^{-\frac{1}{2}} \exp\left(-\frac{1}{2\sigma_\mu^2} (\mu - \mu_0)^2\right)$$

$$\pi(\sigma_\theta^2 | \theta_1, \dots, \theta_k, \mu, \sigma_e^2) = \frac{b_1^{a_1}}{\Gamma(a_1)} (\sigma_\theta^2)^{-a_1-1} \exp\left(-\frac{b_1}{\sigma_\theta^2}\right)$$

$$\pi(\sigma_e^2 | \theta_1, \dots, \theta_k, \mu, \sigma_\theta^2) = \frac{b_2^{a_2}}{\Gamma(a_2)} (\sigma_e^2)^{-a_2-1} \exp\left(-\frac{b_2}{\sigma_e^2}\right)$$

لذا چگالی پیشین و پسین پارامترهای مدل عبارتست از :

$$\pi(\theta_1, \dots, \theta_k, \mu, \sigma_\theta^2, \sigma_e^2) = (2\pi\sigma_\theta^2)^{-\frac{k}{2}} (2\pi\sigma_e^2)^{-\frac{1}{2}} \frac{b_1^{a_1}}{\Gamma(a_1)} \frac{b_2^{a_2}}{\Gamma(a_2)}$$

$$\exp\left(-\frac{1}{2\sigma_\theta^2} \sum_{i=1}^k (\theta_i - \mu)^2 - \frac{1}{2\sigma_e^2} (\mu - \mu_0)^2 - \frac{b_1}{\sigma_\theta^2} - \frac{b_2}{\sigma_e^2}\right)$$

$$\pi(\theta_1, \dots, \theta_k, \mu, \sigma_\theta^2, \sigma_e^2 | \mathbf{y}) \propto f(\mathbf{y} | \theta_1, \dots, \theta_k, \mu, \sigma_\theta^2, \sigma_e^2)$$

$$f(\mathbf{y} | \theta_1, \dots, \theta_k, \mu, \sigma_\theta^2, \sigma_e^2) \pi(\theta_1, \dots, \theta_k, \mu, \sigma_\theta^2, \sigma_e^2)$$

با توجه به فرم پیچیده‌ی توزیع پسین، برای محاسبات عددی استنباط بیز نمی‌توان از آن با روش‌های معمول نمونه‌ای استخراج کرد. اما محاسبات ساده نشان می‌دهد که :

$$\begin{aligned}\theta_i | \mu, \sigma_\theta^2, \sigma_e^2, \mathbf{y} &\sim N\left(\frac{m\sigma_\theta^2 \bar{y}_i + \sigma_e^2 \mu}{m\sigma_\theta^2 + \sigma_e^2}, \frac{\sigma_\theta^2 \sigma_e^2}{m\sigma_\theta^2 + \sigma_e^2}\right) \\ \mu | \theta_1, \dots, \theta_k, \sigma_\theta^2, \sigma_e^2, \mathbf{y} &\sim N\left(\frac{\sigma_\theta^2 \mu_0 + \sigma_e^2 \sum_{i=1}^k \theta_i}{\sigma_\theta^2 + k\sigma_e^2}, \frac{\sigma_\theta^2 \sigma_e^2}{\sigma_\theta^2 + k\sigma_e^2}\right) \\ \sigma_\theta^2 | \theta_1, \dots, \theta_k, \mu, \sigma_e^2, \mathbf{y} &\sim IG(a_2 + \frac{k}{2}, b_2 + \frac{1}{2} \sum_{i=1}^k (\theta_i - \mu)^2) \\ \sigma_e^2 | \theta_1, \dots, \theta_k, \mu, \sigma_\theta^2, \mathbf{y} &\sim IG(a_1 + \frac{k}{2}, b_1 + \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^m (y_{ij} - \theta_i)^2)\end{aligned}$$

وچون می‌توان از توزیعهای نرمال و معکوس گاما نمونه گرفت، لذا آگوریتم گیبز را برای تولید یک نمونه‌ی تقریبی از توزیع پسین بکار می‌بریم. با انتخاب مقدار اولیه‌ی $(\theta_1^{(0)}, \dots, \theta_k^{(0)}, \mu^{(0)}, (\sigma_\theta^2)^{(0)}, (\sigma_e^2)^{(0)})$ ، الگریتم گیبز با داشتن $(\theta_1^{(t+1)}, \dots, \theta_k^{(t+1)}, \mu^{(t+1)}, (\sigma_\theta^2)^{(t+1)}, (\sigma_e^2)^{(t+1)})$ ؛ $(\theta_1^{(t)}, \dots, \theta_k^{(t)}, \mu^{(t)}, (\sigma_\theta^2)^{(t)}, (\sigma_e^2)^{(t)})$ را بصورت زیر تولید می‌کند :

(۱) یک ترتیب از مولفه‌های $(\theta_1^{(t)}, \dots, \theta_k^{(t)}, \mu^{(t)}, (\sigma_\theta^2)^{(t)}, (\sigma_e^2)^{(t)})$ انتخاب می‌کنیم.

(۲) با توجه به ترتیب انتخاب شده در گام یک، θ_i^* را از توزیع

$$N\left(\frac{m(\sigma_\theta^2)^{(t)} \bar{y}_i + (\sigma_e^2)^{(t)} \mu^{(t)}}{m(\sigma_\theta^2)^{(t)} + (\sigma_e^2)^{(t)}}, \frac{(\sigma_\theta^2)^{(t)} (\sigma_e^2)^{(t)}}{m(\sigma_\theta^2)^{(t)} + (\sigma_e^2)^{(t)}}\right)$$

را از توزیع μ^* می‌تولید.

$$N\left(\frac{(\sigma_\theta^2)^{(t)} \mu_0 + \sigma_e^2 \sum_{i=1}^k \theta_i^{(t)}}{(\sigma_\theta^2)^{(t)} + k\sigma_e^2}, \frac{(\sigma_\theta^2)^{(t)} \sigma_e^2}{(\sigma_\theta^2)^{(t)} + k\sigma_e^2}\right)$$

و $(\sigma_e^2)^*$ را از توزیع

$$IG(a_2 + \frac{k}{2}, b_2 + \frac{1}{2} \sum_{i=1}^k (\theta_i^{(t)} - \mu^{(t)})^2)$$

و $(\sigma_e^2)^*$ را از توزیع

$$IG(a_1 + \frac{k}{2}, b_1 + \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^m (y_{ij} - \theta_i^{(t)})^2)$$

تولید می‌کنیم.

(۳) قرار می‌دهیم :

$$(\theta_1^{(t+1)}, \dots, \theta_k^{(t+1)}, \mu^{(t+1)}, (\sigma_\theta^2)^{(t+1)}, (\sigma_e^2)^{(t+1)}) = (\theta_1^*, \dots, \theta_k^*, \mu^*, (\sigma_\theta^2)^*, (\sigma_e^2)^*)$$

«نابع»

- [۱] *Computational statistics; G.Givens, J.Hoeting; Wiley(۲۰۰۵)*
- [۲] *MCMC Methodology; Brani Vidakovic; Lecture*
- [۳] *Numerical Methods Of Statistics; J.F.Monahan; Cambridge University Press(۲۰۰۱)*
- [۴] *Simulation; S.Ross; ۳rd Edition; (۱۹۹۹)*
- [۵] *A Short Course On Computational Statistics; Shojaedin Chenouri; Shahid Beheshti University(Spring ۲۰۰۵)*
- [۶] *General State Space Markov Chains And MCMC Algorithms; G.O.Roberts, J.S.Rosenthal; Probability Surveys(۲۰۰۴)*
- [۷] نخستین درس در فرآیندهای تصادفی؛ ترجمه‌ی : علی اکبر عالم زاده، عین الله پاشا.
- [۸] روش مونت کارلو؛ ترجمه‌ی : یدالله درج.