

# استنباط زنجیره‌های مارکف: مروری بر مقایسه الگو، برآورد بیزی و نرخ آنتروپی

معصومه اصلی بیگی<sup>۱</sup>، حمید پزشک<sup>۲</sup>

چکیده:

زنجیره‌های مارکف در زمینه‌های مختلف علمی برای تحلیل داده‌های دنباله‌ای به کار می‌روند. در این نوشته با استفاده از مقاله [۳] نشان می‌دهیم چگونه می‌توان زنجیره‌های مارکف از مرتبه دلخواه  $k$  را برای داده‌های متناهی با استفاده از روش‌های بیزی در برآورد پارامترها و انتخاب مرتبه الگو مورد استفاده قرار داد. با به کارگیری روش‌هایی از نظریه اطلاع، آنتروپی نسبی و نرخ آنتروپی را محاسبه می‌کنیم. در نهایت با ارائه یک مثال، مراحل بالا را که مراحل مختلف استنباط یک زنجیر مارکف متناهی است به صورت شهودی بیان می‌کنیم.

واژه‌های کلیدی: استنباط بیزی، توزیع پیشین، توزیع پسین، الگوی مارکف، نرخ آنتروپی.

## ۱ مقدمه

می‌کند [۱]. هر سه مرحله نامبرده در ارتباط با یکدیگر هستند و به کارگیری هر یک از آنها بدون توجه به سایر مراحل سبب به وجود آمدن تفسیرهای گمراه کننده می‌شود. ترکیب استنباط پارامترهای الگو، مقایسه الگو و برآورد نرخ آنتروپی [۲ و ۴] به عنوان سومین مرحله، سبب درک بهتر الگوهای زنجیر مارکف می‌شود و این مسئله در مورد منبع تولید داده‌هایی که خارج از کلاس الگوی زنجیر مارکف قرار می‌گیرد نیز درست است. با مقایسه الگو، ساختار منبع داده‌ها معین می‌شود و با برآورد نمودن نرخ آنتروپی توصیفی از تصادفی بودن منبع داده‌ها فراهم می‌شود. این

در بسیاری از سیستم‌های طبیعی، داده‌ها از دنباله‌ای از حروف از یک الفبای متناهی تشکیل شده‌اند. برای مثال، تحلیل اطلاع دنباله در بیویلیمر از این دسته است [۵]. استنباط آماری الگوهایی که در این نوع سیستم‌ها دیده می‌شود از دیرباز مورد توجه محققان بوده است. این الگوها در کلاس الگوهای زنجیرهای مارکف قرار می‌گیرند [۲ و ۴]. در اغلب مطالعات انجام شده برای تجزیه و تحلیل این الگوها، مراحل استنباط آماری به بخش‌های زیر تقسیم می‌شود: ۱. برآورد نقطه‌ای پارامترهای الگو ۲. مقایسه الگو ۳. برآورد توابعی از پارامترهای الگو.

استنباط بیزی همه این مراحل را به هم مرتبط

<sup>۱</sup> کارشناس ارشد آمار، بانک مرکزی  
<sup>۲</sup> استاد گروه آمار دانشگاه تهران

### ۳ زنجیرهای مارکف

اولین قدم در استنباط مشخص کردن فرضیات الگو است. فرض می‌کنیم که یک مجموعه داده  $D$  از طول  $N$ ، نقطه شروع استنباط را نشان دهد.  $D$  از نمادهای  $s_t$  که متعلق به یک الفبای متناهی  $A$  است، انتخاب می‌شود:

$$D = s_0 s_1 \dots s_{N-1} \quad s_t \in A \quad (2)$$

نماد  $s_t^{\leftarrow k}$  برای نشان دادن یک دنباله با طول  $k$  از حروف که در موقعیت  $t$  پایان یافته است معرفی می‌شود: مثلاً  $s_4^{\leftarrow 2} = s_3 s_4$ . کلاس الگوی زنجیر مارکف مرتبه  $k$  ام، یک حافظه متناهی و همگن در منبع داده‌ها متصور می‌شود و می‌تواند به صورت زیر نوشته شود:

$$p(D) = p(s_{k-1}^{\leftarrow k}) \prod_{t=k-1}^{N-1} p(s_{t+1} | s_t^{\leftarrow k}) \quad (3)$$

پارامترهای زنجیر مارکف مرتبه  $k$  ام عبارتند از:

$$\theta_k = \{p(s | s^{\leftarrow k}) : s \in A, s^{\leftarrow k} \in A^k\} \quad (4)$$

برای هر کلمه  $s^{\leftarrow k}$  باید شرط  $\sum_{s \in A} p(s | s^{\leftarrow k}) = 1$  برقرار باشد.

### ۱.۳ درست‌نمایی

با داده‌های  $D = s_0 s_1 \dots s_{N-1}$ ، درست‌نمایی می‌تواند با استفاده از خاصیت مارکف معادله (۳) نوشته شود:

$$P(D | \theta_k, M_k) = \prod_{s \in A} \prod_{s^{\leftarrow k} \in A^k} [p(s | s^{\leftarrow k})]^{n(s^{\leftarrow k}, s)} \quad (5)$$

نوشته که اقتباسی از [۳] است در صدد است نوع متفاوتی از استنباط آماری را دوره کند. این استنباط برای دنباله‌ای متناهی از پیشامدها مورد توجه محققین است.

### ۲ استنباط پارامترهای الگو

در مرحله اول از استنباط بیزی یک ارتباط سیستماتیک بین داده‌ها  $D$ ، الگوی انتخابی  $M$  و بردار پارامترهای الگو  $\theta$  معرفی می‌شود. موضوع جالب توجه در استنباط پارامترهای الگو، چگالی احتمال پسین  $P(\theta | D, M)$  است که احتمال پارامترهای الگو به شرط داده‌های مشاهده شده و الگوی انتخابی است. برای یافتن پسین در ابتدا، چگالی توأم  $P(\theta, D | M)$  داده‌ها و پارامترهای الگو به شرط دانستن اینکه الگوی  $M$  انتخاب شده است، مورد بررسی قرار می‌گیرد. چگالی پسین بنابر قضیه بیز به دست می‌آید:

$$P(\theta | D, M) = \frac{P(D | \theta, M) P(\theta | M)}{P(D | M)} \quad (1)$$

چگالی پیشین  $P(\theta | M)$  توزیع روی پارامترهای الگو را تعیین می‌کند. درست‌نمایی  $P(D | \theta, M)$  احتمال داده‌ها را به شرط پارامترهای الگو و الگوی داده شده شرح می‌دهد. در خاتمه سندیت  $P(D | M)$  احتمال داده‌ها به شرط الگوی داده شده است. در بخش بعدی هر یک از این کمیت‌ها برای به دست آوردن عبارتی صریح برای چگالی پیشین شرح داده می‌شود.

$\delta(t) = 0$  با داشتن این فرم تابعی حداقل دو راه برای تفسیر آنچه که پیشین درباره پارامترهای زنجیر مارکف می‌گوید، وجود دارد: ۱. روش آمار کلاسیک: در این روش میانگین و واریانس  $p(s|s^{+k})$  را نسبت به توزیع پیشین بدست می‌آوریم. ۲. روش آمار بیزی: توزیع حاشیه‌ای را برای پارامتر هر الگو بررسی می‌کنیم. برای توزیع دیریکله توزیع حاشیه‌ای برای هر پارامتر توزیع بتا خواهد بود.

فرض عمومی در استنباط الگو این است که بپذیریم همه پارامترها دارای توزیع پیشین یکنواخت هستند، یعنی:  $\alpha(s^{+k}s) = 1 \quad \forall s \in A, s^{+k} \in A^k$

### ۳.۳ سندیت

با درست‌نمایی داده شده و پیشین انتخاب شده، مدرک  $P(D|M_k)$  جمله نرمال کننده در قضیه بیز است.

$$P(D | M_k) = \int P(D|\theta_k, M_k)P(\theta_k|M_k)d\theta_k \quad (۷)$$

با به کارگیری معادله فوق در درست‌نمایی معادله (۵) و پیشین معادله (۶)، سندیت به دست می‌آید:

$$P(D | M_k) = \prod_{s^{+k} \in A^k} \left\{ \frac{\Gamma[\alpha(s^{+k})]}{\prod_{s \in A} \Gamma[\alpha(s^{+k}s)]} \right. \\ \left. \times \frac{\prod_{s \in A} \Gamma[n(s^{+k}s) + \alpha(s^{+k}s)]}{\Gamma[n(s^{+k}) + \alpha(s^{+k})]} \right\} \quad (۸)$$

که در آن  $n(s^{+k}s)$  تعداد زمان‌هایی است که کلمه  $s^{+k}s$  در نمونه  $D$  رخ داده است. برای استفاده‌های بعدی، نماد  $n(s^{+k}) = \sum_{s \in A} n(s^{+k}s)$  برای تعداد زمان‌هایی که کلمه  $s^{+k}$  مشاهده شده است، معرفی می‌شود. توجه می‌کنیم که معادله (۵) روی دنباله آغازین  $s_{k-1}^{+k} = s_0 s_1 \dots s_{k-1}$  شرطی شده است.

### ۲.۳ چگالی پیشین

چگالی پیشین  $P(\theta_k|M_k)$  برای مشخص کردن فرضیات در مورد الگو، قبل از اینکه داده‌ها در دسترس باشند، استفاده می‌شود. فرم دقیق پیشین به وسیله مفروضات ما از ابرپارامترها  $\alpha(s^{+k}s)$  برای پیشین معین می‌شود. در زنجیر مارکف مرتبه  $k$  ام یک ابرپارامتر برای هر کلمه  $s^{+k}$  به ازای الفبای تحت بررسی داده شده وجود دارد. پیشین مزدوج برای استنباط زنجیر مارکف حاصلضربی از توزیع‌های دیریکله است که در آن برای هر کلمه  $s^{+k}$  چگالی دیریکله‌ای در نظر گرفته شده است [۸]. بنابراین:

$$P(\theta_k|M_k) = \prod_{s^{+k} \in A^k} \left\{ \frac{\Gamma[\alpha(s^{+k})]}{\prod_{s \in A} \Gamma[\alpha(s^{+k}s)]} \right. \\ \times \delta \left( 1 - \sum_{s \in A} p(s|s^{+k}) \right) \\ \left. \times \prod_{s \in A} p(s|s^{+k})^{\alpha(s^{+k}s)-1} \right\} \quad (۶)$$

که در آن  $\alpha(s^{+k}) = \sum_{s \in A} \alpha(s^{+k}s)$  و تابع  $\delta$  عامل نرمال ساز برای پارامترهای الگو است، یعنی به ازای  $t = 0$ ،  $\delta(t) = 1$  و به ازای  $t \neq 0$

### ۴.۳ چگالی پسین

با استفاده از قضیهٔ بیز معادله (۱) و نتایج قبلی، توزیع پسین روی پارامترهای زنجیر مارکف مرتبه  $k$  بدست می‌آید:

$$P(\theta_k | D, M_k) = \prod_{s^{\leftarrow k} \in A^k} \left\{ \frac{\Gamma[n(s^{\leftarrow k}) + \alpha(s^{\leftarrow k})]}{\prod_{s \in A} \Gamma[n(s^{\leftarrow k}s) + \alpha(s^{\leftarrow k}s)]} \times \delta \left( 1 - \sum_{s \in A} p(s | s^{\leftarrow k}) \right) \times \prod_{s \in A} p(s | s^{\leftarrow k})^{n(s^{\leftarrow k}s) + \alpha(s^{\leftarrow k}s) - 1} \right\} \quad (9)$$

چگالی پسین توزیع دیریکله با پارامترهای تغییر یافته نسبت به توزیع پیشین است که نتیجهٔ انتخاب پیشین مزدوج است. همانند چگالی پیشین، دو راه برای فهمیدن آنچه که چگالی پسین درباره تغییرات در پارامترهای زنجیر مارکف برآورد شده به ما می‌گوید، وجود دارد. اولی، استفاده از برآورد نقطه‌ای است و از میانگین و واریانس  $p(s | s^{\leftarrow k})$  نسبت به پسین به دست می‌آید. دوم بررسی توزیع حاشیه‌ای برای پارامتر هر الگو است. برای توزیع دیریکله، توزیع حاشیه‌ای برای هر پارامتر توزیع بتا خواهد بود.

### ۴ مقایسه الگو

مجموعه الگوهای مرتبه  $k$ ،  $M = \{M_k\}_{k_{min}^{k_{max}}}$  را در نظر بگیرید. احتمال توأم  $P(M_k, D | M)$  از یک الگوی بخصوص  $M_k \in M$  و داده‌های  $D$  بنا بر

قضیه بیز به صورت زیر است:

$$P(M_k | D, M) = \frac{P(D | M_k, M) P(M_k | M)}{P(D | M)} \quad (10)$$

که در آن مخرج یک مجموع است و بنا بر قانون کل احتمال با مجموع گیری روی الگوهای  $M_k \in M$  به دست می‌آید. احتمال یک الگوی ویژه در مجموعه تحت بررسی با دو مؤلفه داده می‌شود: مدرک  $P(D | M_k, M)$  که از معادله (۸) مشخص می‌شود و پیشین روی رتبه الگو  $P(M_k | M)$ .

دو پیشین متداول در مقایسه الگو عبارتند از: ۱. پیشین یکنواخت که در آن همه الگوها به طور مساوی محتمل هستند. ۲. الگوهایی که باید برای تعداد پارامترهای آزاد استفاده شده برای برآورد داده‌ها جریمه در نظر بگیرند. در اولین مورد،  $P(M_k | M) = \frac{1}{|M|}$  برای همه مراتب  $k$  یکسان است. احتمال الگوهایی که از این پیشین استفاده می‌کنند، برابر است با:

$$P(M_k | D, M) = \frac{P(D | M_k, M)}{\sum_{M'_k \in M} P(D | M'_k, M)} \quad (11)$$

در مورد دوم که یک جریمه برای تعداد پارامترهای الگو در نظر گرفته می‌شود با استفاده از فرضی که به الگوهای با پارامتر زیاد جریمه بیشتر تعلق می‌دهد، انتخاب الگو را بر اساس رابطه زیر انجام می‌دهیم:

$$P(M_k | D, M) = \frac{P(D | M_k, M) \exp(-|M_k|)}{\sum_{M'_k \in M} P(D | M'_k, M) \exp(-|M'_k|)} \quad (12)$$

که در آن  $|M_k|$  تعداد پارامترهای آزاد الگو است. برای زنجیر مارکف مرتبه  $k$  ام، تعداد پارامترهای

آزاد برابر است با

$$|M_k| = |A|^k (|A| - 1) \quad (۱۳)$$

که در آن  $|A|$  اندازه  $A$  است.

آنتروپی‌های مورد استفاده در معادله (۱۵) به صورت زیر تعریف می‌شوند.

$$H[q(s^{\leftarrow k})] = - \sum_{s^{\leftarrow k}} q(s^{\leftarrow k}) \log_2 q(s^{\leftarrow k}) \quad (۱۶)$$

$$H[q(s^{\leftarrow k})q(s|s^{\leftarrow k})] = \quad (۱۷)$$

$$- \sum_{s^{\leftarrow k}, s} q(s^{\leftarrow k})q(s|s^{\leftarrow k})$$

$$\times \log_2 q(s^{\leftarrow k})q(s|s^{\leftarrow k})$$

## ۵ نظریه اطلاع، مکانیک آماری و نرخ‌های آنتروپی

ارتباط بین استنباط و نظریه اطلاع [۶] با در نظر گرفتن معادله پیشین (۶) و معادله درستیابی (۵) آغاز می‌شود.

$$P(\theta_k|M_k)P(D|\theta_k, M_k) = P(D, \theta_k|M_k) \quad (۱۴)$$

سندیت از رابطه  $Z = P(D|M_k) = \int P(D, \theta_k|M_k)d\theta_k$  به دست می‌آید. با تعریف

$$\beta_k = \sum_{s^{\leftarrow k}, s} [n(s^{\leftarrow k} s) + \alpha(s^{\leftarrow k} s)]$$

و توزیع  $Q$  به صورت

$$Q = \left\{ \begin{aligned} q(s^{\leftarrow k}) &= \frac{n(s^{\leftarrow k}) + \alpha(s^{\leftarrow k})}{\beta_k}, \\ q(s|s^{\leftarrow k}) &= \frac{n(s^{\leftarrow k} s) + \alpha(s^{\leftarrow k} s)}{n(s^{\leftarrow k}) + \alpha(s^{\leftarrow k})} \end{aligned} \right\}$$

می‌توان امید و واریانس  $E(Q, P)$  را با مشتق‌گیری از لگاریتم تابع  $Z$  نسبت به  $\beta_k$  به دست آورد. که در آن  $P$  توزیع پارامترهای واقعی است:  $P = \{p(s^{\leftarrow k}), p(s|s^{\leftarrow k})\}$  فرم مجانبی امید ریاضی  $E(Q, P)$  به صورت زیر خواهد بود.

$$E_{post}[E(Q, P)] = H[q(s^{\leftarrow k})q(s|s^{\leftarrow k})] \quad (۱۵)$$

$$- H[q(s^{\leftarrow k})] + \frac{1}{2\beta_k} |A|^k (|A| - 1) + O\left(\frac{1}{\beta_k^2}\right)$$

با توجه به معادله (۱۵) دو جمله اول نرخ آنتروپی را می‌سازند  $h_\mu[Q] = H[q(s^{\leftarrow k})q(s|s^{\leftarrow k})] - H[q(s^{\leftarrow k})]$  و جمله آخر با آنتروپی نسبی شرطی بین توزیع میانگین پسین  $Q$  و توزیع واقعی  $P$  یکی است. از این رو می‌توان برای برآورد نرخ آنتروپی همگرایی  $E_{post}[E(Q, P)]$  به نرخ آنتروپی  $h_\mu[Q]$  را مورد بررسی قرارداد (برای جزئیات بیشتر به [۳] مراجعه کنید).

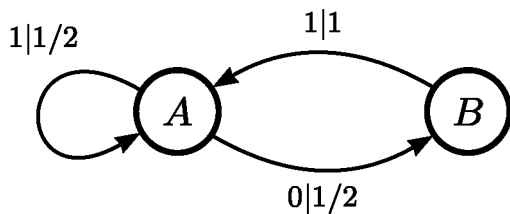
## ۶ مثال‌ها

در این قسمت نظریه‌های بیان شده در قسمت‌های قبل را در قالب یک مثال به صورت کاربردی بیان می‌کنیم. منبع داده‌های این مثال به فرآیند میانگین طلایی معروف است. این فرآیند یک زنجیر مارکف مرتبه اول است [۳]. فرآیند میانگین طلایی می‌تواند به صورت الگویی با دو وضعیت داخلی - آنها را  $A$  و  $B$  بخوانید - ثبت شود. هدف این قسمت، بررسی سه گام اصلی در استنباط برای تجزیه و تحلیل است. در ابتدا استنباط از یک

آنتروپی برای مثال‌ها، از فرمول

$$h_{\mu} = - \sum_{v \in \{A, B\}} p(v) \sum_{s \in A} p(s | v) \log_2 p(s | v)$$

استفاده می‌کنیم. در این فرمول،  $p(s | v) = T_v^{(s)}$  احتمال حرف  $s$  به شرط وضعیت داده شده  $v$  است و  $p(v)$  احتمال مجانبی وضعیت  $v$  است که می‌تواند به صورتی که در بالا بحث شد، محاسبه شود.



شکل ۱. زنجیر مارکف برای فرآیند میانگین طلایی. یال‌ها با نماد خروجی و احتمال انتقال برچسب زده شده‌اند.

فرآیند میانگین طلایی می‌تواند به وسیله یک زنجیر مارکف مرتبه اول روی یک الفبای دودویی نمایش داده شود که در آن هرگز دو صفر متوالی دیده نمی‌شود. ماتریس‌های انتقال برچسب دار تعریف شده برای این منبع داده به صورت زیر هستند.

$$T^{(0)} = \begin{pmatrix} 0 & \frac{1}{2} \\ 0 & 0 \end{pmatrix}, \quad T^{(1)} = \begin{pmatrix} \frac{1}{2} & 0 \\ 1 & 0 \end{pmatrix}$$

شکل ۱ یک نمایش گرافیکی از زنجیر مارکف پنهان منطبق با این فرآیند را نمایش می‌دهد. این شکل وجود یک رابطه ساده بین وضعیت‌های  $A$  و  $B$  و نمادهای خروجی  $0$  و  $1$  را نشان می‌دهد. مشاهده  $0$  نشان دهنده یک انتقال به وضعیت  $B$

زنجیر مارکف مرتبه اول را برای توصیف برآورد پارامترهای الگو با عدم قطعیت، بررسی می‌کنیم. سپس، مقایسه الگو را برای مرتبه‌های مختلف  $k$  بررسی می‌کنیم. این کار سبب می‌شود ساختار منبع داده‌ها را کشف کنیم. در نهایت، نرخ آنتروپی این منبع داده را برآورد می‌کنیم. برای این کار، از فرم شناخته شده منابع سود می‌بریم که هر یک به وسیله یک ماتریس انتقال  $T$  که انتقال بین وضعیت‌های  $A$  و  $B$  را نشان می‌دهد، شرح داده می‌شود.

$$T = \begin{pmatrix} P(A | A) & P(B | A) \\ P(A | B) & P(B | B) \end{pmatrix}$$

وضعیت ویژه  $\vec{\pi} = (P(A), P(B))$  توزیع مجانبی روی وضعیت‌های داخلی را شرح می‌دهد. ماتریس  $T$  می‌تواند به ماتریس‌های برچسب دار  $T^{(s)}$  که شامل عناصری از  $T$  هستند که خروجی آنها نماد  $s$  است تقسیم شود. برای منابع داده‌های دودویی داریم:  $T = T^{(0)} + T^{(1)}$ . با استفاده از این ماتریس‌ها، متوسط احتمال کلمه‌ها را برای هر فرآیند دلخواه برآورد می‌کنیم. برای مثال، احتمال کلمه  $01$  با استفاده از رابطه  $p(01) = \vec{\pi} T^{(0)} T^{(1)} \vec{\eta}$  که در آن  $\vec{\eta}$  یک بردار ستونی است که همه عناصرش  $1$  هستند. برای داده‌های با اندازه  $N$ ، متوسط تعداد دفعاتی که یک کلمه با طول  $k+1$  مشاهده می‌شود را به صورت مقابل برآورد می‌کنیم:  $n(s^{\leftarrow k}s) = (N - k)p(s^{\leftarrow k}s)$

است و مشاهده ۱ دلالت بر یک انتقال به وضعیت  $A$  دارد و لذا این فرآیند یک زنجیر مارکف روی  $\circ$  ها و ۱ ها ایجاد می‌کند.  $p(A)$  برابر است با:  $p(B) = p(\circ) = \pi T^{(\circ)} \eta = \frac{1}{3}$  و  $p(A) = p(1) = \pi T^{(1)} \eta = \frac{2}{3}$  با: فرآیند داریم  $\pi = [p(A), p(B)] = (\frac{2}{3}, \frac{1}{3})$  است. با این بردار و ماتریس‌های انتقال برچسب دار تعریف شده در بالا، متوسط تعداد دفعات رخداد هر کلمه خواسته شده می‌تواند مطابق آنچه در بالا شرح داده شد به دست آید.

## ۱.۶ برآورد پارامترها

به منظور استنباط پارامترهای زنجیر مارکف برای فرآیند میانگین طولی از متوسط تعداد دفعات تکرار هر کلمه در اندازه‌های متنوع  $N$  یعنی  $n(s^{\leftarrow k} s)$  استفاده می‌کنیم. برای این کار روش دوم، یعنی بررسی توزیع حاشیه‌ای برای پارامترهای الگورا به کار می‌گیریم. با توجه به معادله (۹) چگالی پسین حاشیه‌ای برای توزیع دیریکله توزیع

بتا است و برای پارامترهای  $p(\circ|1)$  و  $p(1|\circ)$  به ترتیب به صورت زیر است.

$$P(p(\circ|1)) = \frac{\Gamma[n(1) + \alpha(1)]}{\Gamma[n(1|\circ) + \alpha(1|\circ)]} \times \frac{p(\circ|1)^{n(1|\circ) + \alpha(1|\circ) - 1}}{\Gamma[n(1) + \alpha(1) - n(1|\circ) - \alpha(1|\circ)]} \times (1 - p(\circ|1))^{n(1) + \alpha(1) - n(1|\circ) - \alpha(1|\circ) - 1}$$

$$P(p(1|\circ)) = \frac{\Gamma[n(\circ) + \alpha(\circ)]}{\Gamma[n(\circ|1) + \alpha(\circ|1)]} \times \frac{p(1|\circ)^{n(\circ|1) + \alpha(\circ|1) - 1}}{\Gamma[n(\circ) + \alpha(\circ) - n(\circ|1) - \alpha(\circ|1)]} \times (1 - p(1|\circ))^{n(\circ) + \alpha(\circ) - n(\circ|1) - \alpha(\circ|1) - 1}$$

برای نمونه‌های با اندازه ۴۰۰ و ۲۰۰ و ۱۰۰ و  $N=50$ ، چگالی پسین حاشیه‌ای را برای پارامترهای  $p(\circ|1)$  و  $p(1|\circ)$  بدست می‌آوریم. برای این کار ابتدا مقادیر  $n(s^{\leftarrow k} s)$  را از رابطه  $n(s^{\leftarrow k} s) = (N - K)p(s^{\leftarrow k} s)$  به ازای  $k = 1$  می‌یابیم. برای مثال به ازای  $N=50$  مقادیر فوق به صورت جدول (۱) است:

جدول ۱. مقادیر  $p(s^{\leftarrow k} s)$  و  $n(s^{\leftarrow k} s)$  به ازای  $N=50$ .

$p(\circ \circ) = \pi T^{(\circ)} T^{(\circ)} \eta = \circ$	$n(\circ \circ) = \circ$
$p(\circ 1) = \pi T^{(\circ)} T^{(1)} \eta = \frac{1}{3}$	$n(\circ 1) = \frac{49}{3} = 16/33$
$p(1 \circ) = \pi T^{(1)} T^{(\circ)} \eta = \frac{2}{3}$	$n(1 \circ) = \frac{49}{3} = 16/33$
$p(1 1) = \pi T^{(1)} T^{(1)} \eta = \frac{2}{3}$	$n(1 1) = \frac{49}{3} = 16/33$

بتا با پارامترهای (۱) و (۱۷/۳) است:

$$P(p(1|\circ)) = \frac{\Gamma(18/33)}{\Gamma(17/33)\Gamma(1)} \times p(1|\circ)^{16/33} (1 - p(1|\circ))^\circ$$

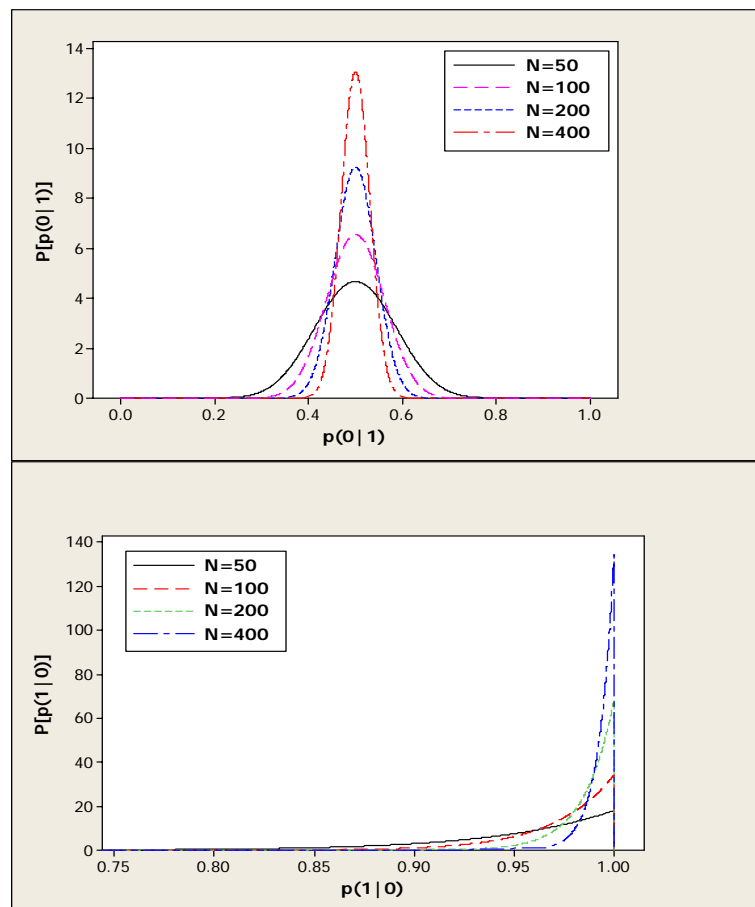
چگالی پسین حاشیه‌ای برای پارامترهای  $p(\circ|1)$  و  $p(1|\circ)$  در شکل ۲ رسم شده است.

چگالی پسین حاشیه‌ای برای پارامتر  $p(\circ|1)$  توزیع

بتا با پارامترهای (۱۷/۳) و (۱۷/۳) است:

$$P(p(\circ|1)) = \frac{\Gamma(34/66)}{\Gamma(17/33)\Gamma(17/33)} \times p(\circ|1)^{16/33} (1 - p(\circ|1))^{16/33}$$

چگالی پسین حاشیه‌ای برای پارامتر  $p(1|\circ)$  توزیع



شکل ۲. نمودار استنباط پارامترهای الگوی  $M_1$  برای فرآیند میانگین طلایی. برای هر نمونه داده به حجم  $N$ ، چگالی پسین حاشیه ای برای پارامترهای مورد علاقه رسم شده است: در نمودار بالا  $p(0|1)$  و در نمودار پایین  $p(1|0)$  رسم شده است. مقادیر واقعی پارامترها  $\frac{1}{4}$ ،  $p(0|1) = 1$ ،  $p(1|0) = 1$  هستند.

واقعی تری از فرآیند استنباط نسبت به برآورد حداکثر درست‌نمایی، با خطای برآورد شده که تنها به وسیله تقریب گوسی از درست‌نمایی به دست می‌آید، ارائه می‌دهد. همان‌طور که شکل ۲ نشان می‌دهد، در حقیقت، تقریب گوسی از عدم قطعیت یک توصیف ناکارآمد از دانسته‌های ما زمانی است که پارامترهای زنجیر مارکف نزدیک حدود بالایی یا پائینی خود در ۰ و ۱ هستند. شاید نتیجه بخش ترین مجموعه از مقادیری که می‌توان به دست

نمودارها توزیع‌های ممکن برای پارامترهای الگو را برای هر  $N$  توصیف می‌کنند. همان‌طور که می‌بینیم با افزایش حجم داده‌ها به مقادیر صحیح  $\frac{1}{4}$  و  $p(0|1) = 1$  و  $p(1|0) = 1$  همگرا می‌شوند. برآوردهای نقطه‌ای با واریانس‌های متناظر می‌تواند برای هر یک از پارامترها به دست آید، ولی این مقادیر خود می‌توانند گمراه کننده باشند. به هر حال، برآورد به دست آمده از به کارگیری میانگین و واریانس توزیع پسین توصیف

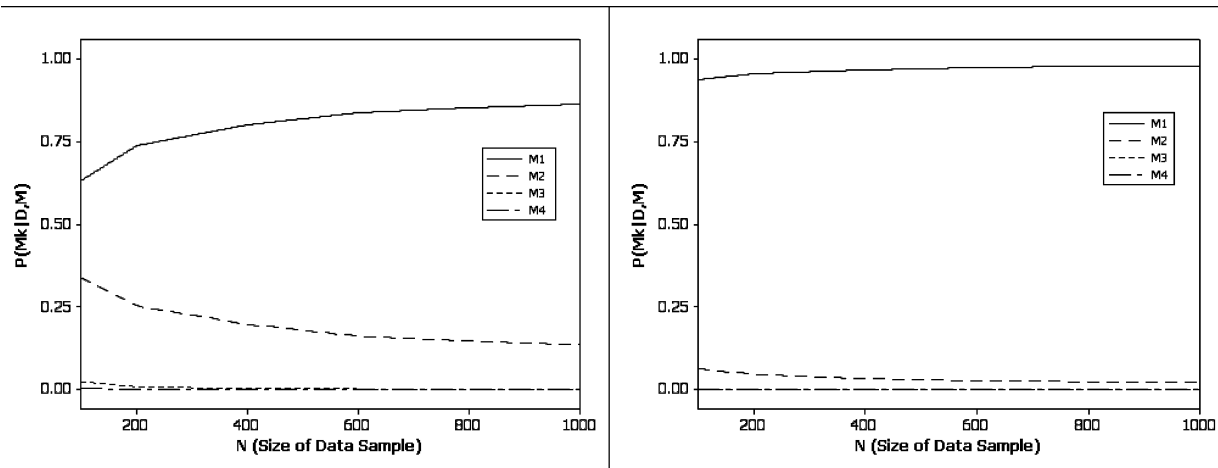


نتیجه انتظار داریم مقایسه الگو، این مرتبه را از میان احتمال‌های مورد بررسی انتخاب نماید. برای شرح این مسئله، مرتبه‌های  $k = 1, \dots, 4$  را بررسی می‌کنیم و مقایسه الگو را با توزیع پیشین یکنواخت روی مرتبه‌ها (معادله ۱۱) و نیز با یک جریمه برای تعداد پارامترهای الگو (معادله ۱۲) انجام می‌دهیم. نتایج مقایسه الگو در شکل ۳ نشان داده شده است.

آورد شامل میانگین توزیع پسین و ناحیه اطمینان باشد، اینها با بیشترین دقت، عدم تقارن در عدم قطعیت پارامترهای الگو را توصیف خواهند کرد.

## ۲.۶ انتخاب رتبه الگو

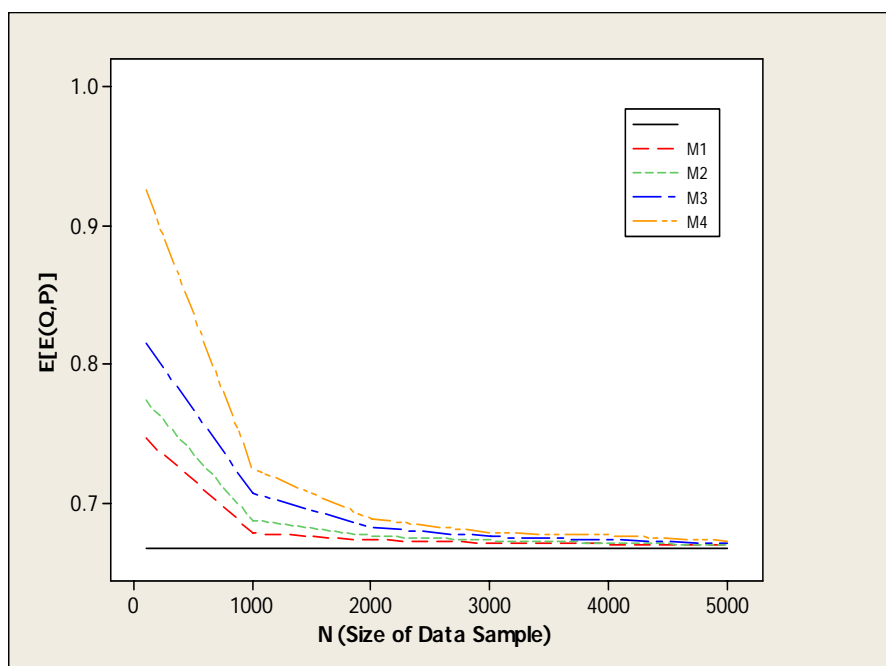
همانطور که بیان کردیم فرآیند میانگین طلایی یک زنجیر مارکف مرتبه اول است ( $k = 1$ ). در



شکل ۳. مقایسه الگو برای زنجیره‌های مارکف از مرتبه  $k = 1, \dots, 4$  با استفاده از مقادیر متوسط فرآیند میانگین طلایی. حجم‌های نمونه از  $N = 100$  تا  $N = 1000$  در گام‌هایی از  $\Delta N = 5$  برای تولید این نمودارها استفاده شده‌اند.

توزیع پیشین با جریمه روی مرتبه  $k$  نشان می‌دهد. این توزیع پیشین بیش برآزش شده را در حجم‌های کوچک داده‌ها برطرف می‌کند. البته این مطلب به معنی بهتر بودن این توزیع پیشین نیست، بلکه در حقیقت، مقایسه الگوی بی‌زی با توزیع پیشین یکنواخت یک کار مؤثر با به کارگیری حجم نمونه نسبتاً کوچک است.

نمودار سمت چپ احتمال هر مرتبه  $k$  را به صورت تابعی از حجم نمونه با به کارگیری توزیع پیشین یکنواخت، نشان می‌دهد. با این توزیع پیشین روی مرتبه‌ها،  $M_1$  با هر اندازه مناسب از داده‌ها انتخاب شده است. هرچند به نظر می‌رسد امکان بیش برآزش برای حجم کم داده‌ها ( $N \leq 100$ ) وجود دارد. نمودار سمت راست احتمال الگو را با



شکل ۴. همگرایی  $E_{post}[E(Q, P)]$  به نرخ واقعی آنروپی  $h_\mu = \frac{2}{3}$  افقی نشان داده شده است (در شکل با خبرای فرآیند میانگین طلایی). همانطور که در معادله (۱۵) شرح داده شد، آنروپی نسبی شرطی  $D[Q||P] \rightarrow 0$  و این همگرایی در همگرایی  $h_\mu[Q]$  به نرخ واقعی آنروپی نتیجه می‌شود.

### ۳.۶ برآورد نرخ آنروپی

زنجیره‌های مارکف مراتب بالاتر است. در ارزیابی مقدار  $D[Q||P] + h_\mu[Q]$  برای نمونه‌های با اندازه متفاوت، انتظار داریم که برآورد میانگین توزیع پسین  $Q$  به توزیع واقعی  $P$  همگرا شود. در نتیجه باید آنروپی نسبی شرطی با افزایش  $N$  به صفر میل کند. برای فرآیند میانگین طلایی، مقدار واقعی نرخ آنروپی برابر با  $h_\mu = \frac{2}{3}$  است. شکل ۴ همگرایی مورد انتظار از میانگین  $E(Q, P)$  به نرخ واقعی آنروپی را نشان می‌دهد. نتیجه‌ای که از مقایسه الگو در بخش قبل گرفتیم، می‌تواند در برآورد نرخ آنروپی نیز مورد استفاده قرار گیرد. همان‌طور که در شکل ۳ دیدیم، دامنه‌هایی از اندازه نمونه  $N$  وجود دارند که در آن‌ها احتمال مرتبه‌های

حال می‌توانیم همگرایی میانگین  $E(Q, P) = D[Q||P] + h_\mu[Q]$  را به نرخ درست آنروپی برای فرآیند میانگین طلایی شرح دهیم. این همگرایی را برای همه مرتبه‌های  $k = 1, \dots, 4$  که در قسمت قبل بحث شد، نشان می‌دهیم. این مثال نشان می‌دهد که همه مرتبه‌های بزرگ‌تر یا مساوی  $k = 1$  به نرخ آنروپی میل می‌کنند. هرچند همگرایی به مقادیر واقعی برای مرتبه‌های بالاتر به خاطر بزرگ بودن مقدار اولیه  $D[Q||P]$  دربرگیرنده داده‌های بیشتری است. این مقدار بزرگ‌تر به دلیل تعداد بزرگ‌تر پارامترها برای

بررسی کردیم. در اغلب روش‌های استنباط، این سه وجه به صورت مجزا ولی در ارتباط با هم به کار می‌روند. در این مقاله نیز آنها را بسیار وابسته به هم یافتیم. برآورد پارامترهای الگو بدون در نظر گرفتن اینکه آیا الگوی صحیح انتخاب شده است یا نه، در بهترین حالت نیز گمراه کننده است. مقایسه الگو به وسیله مقایسه مرتبه‌های متنوع  $k$  در حدود کلاس الگو راهی به سوی حل این مسئله است. در نهایت، برآورد تصادفی بودن با استفاده از نرخ‌های آنتروپی، اطلاع بیشتری درباره رابطه بین ساختار و تصادفی بودن تولید می‌کند. این بینش‌ها علی‌رغم داده‌های خارج از کلاس، توان ترکیب این سه روش را به صورت یک ابزار مؤثر برای بررسی ساختار و تصادفی بودن در رشته‌های متنهای از داده‌های گسسته، شرح می‌دهند.

$k = 1, 2$  هر دو غیر صفر هستند. در اصل، برآورد  $h_\mu$  باید با وزن دار کردن مقادیر به دست آمده برای هر  $k$  به وسیله احتمال مرتبه مطابق با آن  $P(M_k | D, M)$  مشخص شود. همانطور که از شکل ۴ می‌بینیم، برآوردهای نرخ آنتروپی برای  $k = 1, 2$  در این محدوده از  $N$  نیز بسیار مشابه هستند. در نتیجه، این روش اضافی نیایستی تأثیر بزرگی بر برآورد نرخ آنتروپی داشته باشد.

#### ۴.۶ نتیجه گیری

در این مقاله با استفاده از مقاله [۳]، استنباط بیزی از الگوهای زنجیر مارکف مرتبه  $k$  را که شامل برآورد پارامترهای الگو برای یک  $k$  داده شده به همراه مقایسه الگو بین مرتبه‌ها و برآورد تصادفی بودن با استفاده از نرخ‌های آنتروپی بود،

## مراجع

- [1] Thornburg, H. (2000), *Introduction to Bayesian Statistics*, Standford, California: Standford University.
- [2] Anderson, T.W. and Goodman, L.A. (1957), Statistical inference about Markov chains, *The Annals Mathematical Statistics*, 28, 89-110.
- [3] Strelhoff, C.C., Crutchfield, J.P. and Hubler, A.W. (2007), Inferring Markov chains: Bayesian estimation, model comparison, entropy rate, and out-of-class modeling, *Phys. Rev. E.*, **76**, 011106.1-011106.14.
- [4] Chatfield, C. (1973), Statistical inference regarding markov chain models, *Applied Statistics*, 22, 7-20.

- 
- [5] Liu, J.S. and Lawrence, C.E. (1999), Bayesian inference on biopolymer models, *Bioinformatics*, 15, 38-52.
- [6] Cover, T.M. and Thomas, J.A. (1991), *Elements of Information Theory*, Wiley, New York.
- [7] Shannon, C.E. (1948), A mathematical theory of communication, *Bell System Technical Journal*, 27, 379-423.
- [8] MacKay, D.J.C. and Peto, L.C.B. (1994), A hierarchical dirichlet language model, *Nat. Lang. Eng.*, 1, 1-19.