

بستگی احتمالاتی و بسندگی الگوریتمی از دیدگاه نظریه اطلاع

مهدی شمس^۱

تاریخ دریافت: ۱۳۹۵/۱۰/۷

تاریخ پذیرش: ۱۳۹۶/۱/۳۰

چکیده:

با توجه به اهمیت زنجیر مارکوف در نظریه اطلاع، تعریف احتمال شرطی این فرایند تصادفی می‌تواند بر حسب اطلاع متقابل نیز تعریف شود. در این مقاله ارتباط بین مفهوم بسندگی و زنجیر مارکوف از دیدگاه اصول نظریه اطلاع، و همچنین ارتباط بین بسندگی احتمالاتی و بسندگی الگوریتمی مشخص می‌شود.

واژه‌های کلیدی: تابع بسنده، زنجیر مارکوف، اطلاع متقابل، بسندگی احتمالاتی، بسندگی الگوریتمی.

۱ مقدمه

در این مقاله با تعریف تابع بسنده و بررسی مفهوم آن از طریق احتمال شرطی، ارتباط بین مفهوم بسندگی و زنجیر مارکوف بیان می‌شود. همچنین ارتباط بین بسندگی احتمالاتی و بسندگی الگوریتمی تعیین می‌گردد و ثابت می‌شود که اگر یک آماره بسنده الگوریتمی باشد، این آماره، بسنده احتمالاتی نیز هست. بسندگی احتمالاتی در شاخه‌هایی نظیر قابلیت اطمینان [۱۷] و آمار فازی [۱۳] کاربرد دارد. بسندگی الگوریتمی هم در شاخه‌هایی مثل خوشه‌بندی [۴] و یادگیری ماشین [۲۲] مورد استفاده قرار می‌گیرد. در بخش ۳ خواننده باید با مفاهیمی از نظریه کدگذاری آشنا باشد، که برای درک بهتر این بخش، مطالعه [۱۲] و فصل ۱۴ از [۳] توصیه می‌شود.

فضای احتمال (Ω, \mathcal{F}, P) را در نظر بگیرید، که فضای نمونه Ω متناهی است. متغیر تصادفی گسسته X با مقادیر x_1, \dots, x_N را در نظر بگیرید که $P(X = x_k) = p_k$ ، $k = 1, \dots, N$ آنتروپی X به صورت $H(X) = \sum_{k=1}^N p_k \log_2 \frac{1}{p_k}$ تعریف می‌شود [۳]. از این رو $H(X) \geq 0$ و شرط لازم و کافی برای $H(X) = 0$ آن است که $N = 1$ ، یعنی X یک متغیر تصادفی تباهیده باشد. آنتروپی، مقدار اطلاعات دریافت شده از X را بعد از مشاهده مقدار واقعی X نشان می‌دهد. حال اگر f تابعی با دامنه $\{x_1, \dots, x_N\}$ باشد به طوری که برای $i \neq j$ داشته باشیم $f(x_i) \neq f(x_j)$ ، خواهیم داشت $P(f(X) = f(x_k)) = P(X = x_k) = p_k$ و از این رو

امروزه نظریه فرایندهای مارکوف یک بخش اصلی از نظریه احتمال است و کاربردهای گسترده‌ای در بسیاری از علوم دیگر دارد [۲۳]. مارکوف [۱۴] آغازکننده این نظریه در سال ۱۹۰۷ بوده است [۱]. یکی از اولین کاربردهای این تئوری، تحقیق روی حروف صدادار و بی‌صدا در زبان روسی بوده است [۱۶]؛ اگرچه زبان‌شناسان مدرن نشان داده‌اند که در برخی زبان‌ها خصوصیت زنجیر مارکوف برقرار نیست. شاخه نظریه اطلاع در سال ۱۹۴۸ [۱۸] توسط شانون معرفی شد و در این شاخه از زنجیرهای مارکوف به‌عنوان یک ابزار استفاده می‌شود. یکی از کاربردهای نظریه اطلاع زمانی است که احتمال وقوع یک نماد در یک پیام به تعداد متناهی از نمادهای قبلی بستگی داشته باشد. در این حالت می‌توان از منبع اطلاع با حافظه^۲ نام برد و دنباله تولید شده توسط چنین منبعی را یک زنجیر مارکوف در نظر گرفت. در [۷] ارتباط بین اطلاع شرطی و بسندگی زیرمیدان‌ها بررسی شده و ثابت می‌شود اطلاع شرطی برابر با اختلاف بین اطلاعات حاوی میدان و زیرمیدان است. در [۱۰] ارتباط بین بسندگی و اطلاع از این جهت که اطلاعات نمونه برای آماره بسنده نسبت به نمونه اصلی نه کم می‌شود و نه زیاد، مورد بررسی قرار گرفته است. در [۲۴] راجع به ارتباط بین مفاهیم نسبت درست‌نمایی مجانبی و آنتروپی و بسندگی، مطالبی گفته شده است. در [۲۲، ۱۰] آماره بسنده مینیمال الگوریتمی معرفی شده است.

^۱ عضو هیئت علمی گروه آمار، دانشگاه کاشان، ایران

^۲ information source with memory

برقراری تساوی آن است که X_1, \dots, X_N مستقل باشند [۳]. اطلاع متقابل فراهم آمده توسط متغیرهای تصادفی X_1, \dots, X_N به صورت

$$I(X_1, \dots, X_N) = H(X_1, \dots, X_N) - \sum_{k=1}^N H(X_k)$$

تعریف می‌شود. شرط لازم و کافی برای استقلال X_1, \dots, X_N آن است که اطلاع متقابل آنها صفر باشد یا به طور معادل برای هر $k = 1, \dots, N-1$ داشته باشیم $I((X_1, \dots, X_k), X_{k+1}) = 0$. [۳] اگر متغیرهای تصادفی مستقل X_1, \dots, X_N و متغیر تصادفی Z تعداد متناهی از مقادیر مجزایی را اختیار کنند، ثابت می‌شود [۸]:

$$\sum_{k=1}^N I(X_k, Z) \leq I((X_1, \dots, X_k), Z).$$

تعریف ۲.۱. [۳] متغیرهای تصادفی X و Y و Z از یک زنجیر مارکوف می‌آیند و می‌نویسند $X \rightarrow Y \rightarrow Z$ هرگاه

$$p(x, y, z) = p(x)p(y|x)p(z|y).$$

با استفاده از تعریف بالا نتیجه زیر حاصل می‌شود:

نتیجه ۳.۱. [۳] اگر X و Y و Z متغیرهای تصادفی باشند، داریم:

۱. $X \rightarrow Y \rightarrow Z$ اگر و تنها اگر X و Z به شرط Y مستقل باشند.

۲. $X \rightarrow Y \rightarrow Z$ نتیجه می‌دهد $Z \rightarrow Y \rightarrow X$ و بنا بر این می‌توان نوشت $X \leftrightarrow Y \leftrightarrow Z$.

۳. اگر $Z = f(Y)$ ، آن‌گاه $X \rightarrow Y \rightarrow Z$.

گزاره زیر، یکی از نامساوی‌های مهم در نظریه اطلاع با نام نامساوی پردازش داده‌ها است که برای معرفی مفهوم بسندگی از دیدگاه نظریه اطلاع به آن نیاز داریم.

گزاره ۴.۱. (نامساوی پردازش داده‌ها) [۳] اگر $X \rightarrow Y \rightarrow Z$ آن‌گاه $I(X, Y) \geq I(X, Z)$

با استفاده از گزاره ۴.۱ حقایق زیر به دست می‌آیند [۳]:

۱. اگر $Z = f(Y)$ ، آن‌گاه $I(X, Y) \geq I(X, f(Y))$

۲. اگر $X \rightarrow Y \rightarrow Z$ ، آن‌گاه $I(X, Y) \geq I(X, Y|Z)$

$H(f(X)) = H(X)$. در حالت کلی‌تر برای هر تابع f داریم $H(f(X)) \leq H(X)$. از این که $H(X)$ فقط به توزیع $\mathcal{P} = \{p_1, \dots, p_N\}$ مربوط به X بستگی دارد، برخی مواقع آنتروپی به صورت $H(\mathcal{P}) = \sum_{k=1}^N p_k \log_2 \frac{1}{p_k}$ نیز نوشته می‌شود. حال اگر $\mathcal{Q} = \{q_1, \dots, q_N\}$ احتمال‌های مربوط به متغیر تصادفی دیگری مانند Y باشد و $N \geq 2$ ، واگرایی اطلاع توزیع \mathcal{P} از \mathcal{Q} را به صورت $D(\mathcal{P} \parallel \mathcal{Q}) = \sum_{k=1}^N p_k \log_2 \frac{p_k}{q_k}$ تعریف می‌کنند که آن را فاصله کولبک-لایبلا [۱۰] یا آنتروپی نسبی یا آنتروپی متقابل^۴ نیز می‌گویند [۳].

گزاره ۱.۱. [۳] برای دو توزیع دلخواه \mathcal{P} و \mathcal{Q} داریم $D(\mathcal{P} \parallel \mathcal{Q}) \geq 0$ و تساوی برقرار است اگر و تنها اگر برای هر $k = 1, \dots, N$ داشته باشیم $p_k = q_k$

برای متغیر تصادفی گسسته X ، آنتروپی بیشین است اگر و تنها اگر برای $k = 1, \dots, N$ داشته باشیم $p_k = \frac{1}{N}$ و در این حالت $H(X) = \log_2 N$. [۳] اطلاع متقابل^۵ متغیرهای تصادفی X و Y که هر یک تعداد متناهی از مقادیر مجزا را اختیار می‌کند به صورت $I(X, Y) = H(X) + H(Y) - H(X, Y)$ تعریف می‌شود، که $H(X, Y)$ آنتروپی (X, Y) است [۸]. از این رو $I(X, Y) \geq 0$ و شرط لازم و کافی برای برقراری تساوی، استقلال X و Y است [۸]. کمیت $I(X, Y)$ مقدار اطلاع به دست آمده راجع به متغیر تصادفی X بر اساس مشاهدات متغیر تصادفی Y را اندازه می‌گیرد [۳]. اگر \mathcal{P}_k معرف خانواده توزیع‌های شرطی $\{p_{j|k}, j = 1, \dots, N\}$ باشد که در آن $p_{j|k} = P(X = x_j | Y = y_k)$ میانگین آنتروپی شرطی X به شرط Y به صورت $H(X|Y) = \sum_{k=1}^N q_k H(\mathcal{P}_k)$ تعریف می‌شود و می‌توان نشان داد [۳] که

$$H(X|Y) = H(X, Y) - H(Y) \geq 0,$$

$$I(X, Y) = H(X) - H(X|Y),$$

$$I(X, Y) \leq \min(H(X), H(Y)).$$

برای متغیرهای تصادفی $\{X_i : i = 1, \dots, N\}$ که تعداد متناهی از مقادیر مجزا را اختیار می‌کنند و $N > 2$ ، داریم $H(X_1, \dots, X_N) \leq \sum_{k=1}^N H(X_k)$ و شرط لازم و کافی برای

^۳ Kullback-Leibler distance

^۴ cross entropy

^۵ mutual information

فرض کنید $\{P_\theta\}$ خانواده‌ای از توزیع‌ها مربوط به متغیر تصادفی X باشد که مقادیر خود را در مجموعه شمارای \mathcal{X} اختیار می‌کند و Θ فضای پارامتر باشد. تابع $T: \mathcal{X} \rightarrow \mathcal{T}$ را یک آماره از داده‌ها در \mathcal{X} نامند. در این حالت $\theta \rightarrow X \rightarrow T(X)$ و در پی آن با به کارگیری گزاره ۴.۱ داریم:

$$I(\theta, X) \geq I(\theta, T(X)). \quad (1)$$

اکنون این سؤال مطرح می‌شود که آیا می‌توان آماره‌ای انتخاب کرد که همه اطلاعات X راجع به θ را حفظ کند؟ آماره T برای خانواده $\{P_\theta\}$ (یا برای θ) بسنده است هرگاه برای هر $t \in \mathcal{T}$ ، توزیع شرطی $f_\theta(x|t) = P_\theta(X=x|T(X)=t)$ مستقل از θ باشد. به عبارت دیگر، این توزیع شرطی تحت تغییر θ پایا است [۱۱]. این تعریف نشان می‌دهد همه اطلاعاتی که در مشاهده x راجع به پارامتر θ است، در آماره $T(x)$ خلاصه شده است. طبق قضیه تجزیه فیشر-نیمن [۱۱] داریم $f_\theta(x) = f(x|t)f_\theta(t)$ ، یا به‌ازای $s \in \Theta$ و $x \in \mathcal{X}$

$$P(X=x|\theta=s) = P(X=x|T(X)=T(x)) \times P(T(X)=T(x)|\theta=s).$$

از این رو به‌ازای هر $s \in \Theta$ و $x \in \mathcal{X}$

$$\begin{aligned} P(\theta=s, X=x|T(X)=T(x)) &= \frac{P(\theta=s, X=x)}{P(T(X)=T(x))} \\ &= \frac{P(X=x|\theta=s)P(\theta=s)}{P(T(X)=T(x))} \\ &= \frac{P(X=x|T(X)=T(x))P(T(X)=T(x)|\theta=s)P(\theta=s)}{P(T(X)=T(x))} \\ &= \frac{P(X=x, T(X)=T(x))}{P(T(X)=T(x))} \frac{P(T(X)=T(x), \theta=s)}{P(T(X)=T(x))} \\ &= P(\theta=s|T(X)=T(x))P(X=x|T(X)=T(x)), \end{aligned}$$

و از این رو بسندگی $T(X)$ معادل با این حقیقت است که برای هر توزیع θ روی Θ ، θ و X به شرط $T(X)$ از هم مستقل هستند. با استفاده از این حقیقت و با به کارگیری قسمت (۱) و (۲) از نتیجه ۳.۱، $\theta \rightarrow T(X) \rightarrow X$ و در پی آن با استفاده از نامساوی پردازش داده‌ها داریم:

$$I(\theta, T(X)) \geq I(\theta, X). \quad (2)$$

علاوه بر این از استقلال شرطی θ و X به شرط $T(X)$ مشاهده

می‌شود:

$$\begin{aligned} \frac{P(\theta=s, X=x)}{P(T(X)=T(x))} &= \frac{P(\theta=s, T(X)=T(x))}{P(T(X)=T(x))} \frac{P(X=x)}{P(T(X)=T(x))}, \end{aligned}$$

و بالاخره $P(\theta=s|X=x) = P(\theta=s|T(X)=T(x))$ که معادله اخیر این حقیقت را تصدیق می‌کند که اگر به‌جای X از آماره بسنده $T(X)$ استفاده شود اطلاعی راجع به پارامتر θ از دست نمی‌رود.

در پایان با استفاده از نامساوی‌های (۱) و (۲) برای یک آماره بسنده دلخواه T داریم:

$$I(\theta, X) = I(\theta, T(X)). \quad (3)$$

ذکر این نکته ضروری است که بسندگی نمی‌تواند گسترش یابد؛ به این مفهوم که با توجه به رابطه (۱) اطلاع متقابل داده‌ها و مدل به‌هیچ‌وجه نمی‌تواند توسط پردازش داده‌های نمونه افزایش یابد.

به‌طور معادل با استفاده از تعریف بسندگی می‌توان گفت که شرط لازم و کافی برای این که آماره T بسنده باشد آن است که تابع $g: \mathcal{X} \times \mathcal{T} \rightarrow \mathbb{R}$ وجود داشته باشد به‌طوری که به‌ازای هر $\theta \in \Theta$ و $t \in \mathcal{T}$ و $x \in \mathcal{X}$

$$g(x, t) = f_\theta(x|t). \quad (4)$$

برای درک و فهم بیشتر بسندگی از دیدگاه نظریه اطلاع، مثال زیر می‌تواند مفید باشد.

مثال ۵.۱. [۳] فرض کنید $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{ber}(\theta)$ که $\theta \in (0, 1)$ بنا بر این،

$$f_\theta(\mathbf{x}) = \theta^{T(\mathbf{x})}(1-\theta)^{n-T(\mathbf{x})},$$

$$\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X} = \{0, 1\}^n,$$

که در آن $T(\mathbf{x}) = \sum_{i=1}^n x_i$ تعداد ۱ها در \mathbf{x} است. در این حالت،

حقیقت معادل با کمینه کردن تابع لاگرانژ

$$\mathcal{L}[p(t|x)] = I(X, T) - \beta I(Y, T)$$
 نسبت به شرط
 مارکوف است.

متغیر تصادفی T را که در شرایط بالا صدق می‌کند، گره
 اطلاع^۶ بین X و Y گویند. حل مسئله گره اطلاع به معادلات
 زیر منجر می‌شود [۲۰]:

$$p(t|x) = \frac{p(t)}{Z(x, \beta)} \exp(-\beta D(p(y|x) \| p(y|t))),$$

که در آن،

$$p(t) = \sum_x p(t|x)p(x)$$

$$p(y|t) = \sum_x p(y|x)p(x|t)$$

$$Z(x, \beta) = \sum_t p(t) \exp(-\beta D(p(y|x) \| p(y|t))),$$

$$D(p(y|x) \| p(y|t)) = E_{p(y|x)} \left[\log \frac{p(y|x)}{p(y|t)} \right].$$

۲ ارتباط زنجیرهای مارکوف و بسندگی از دیدگاه نظریه اطلاع

عنصر اصلی فرمول‌بندی خاصیت مارکوف بر حسب آنتروپی و
 اطلاع، استفاده از دیدگاه تابع بسنده است. ایده اصلی بسندگی
 یک متغیر تصادفی برای یک متغیر تصادفی دیگر از رابطه (۳)
 گرفته شده است.

تعریف ۱.۲. [۱۶] اگر X و Y متغیرهای تصادفی روی یک
 فضای احتمال (Ω, \mathcal{F}, P) باشند به طوری که $I(X, Y) < \infty$ ، و
 g یک تابع با مقدار حقیقی و بول اندازه‌پذیر باشد، متغیر تصادفی
 $g \circ X = g(X)$ را یک تابع بسنده از متغیر تصادفی X
 برای متغیر تصادفی Y گویند هرگاه $I(g(X), Y) = I(X, Y)$
 به عبارت دیگر $g(X)$ برای Y بسنده است هرگاه $g(X)$ شامل
 همه اطلاعات متغیر تصادفی Y باشد که توسط متغیر تصادفی X
 فراهم می‌آیند.

قضیه زیر، یک شرط لازم و کافی برای بسندگی یک آماره
 برای یک متغیر تصادفی دیگر ارائه می‌کند.

قضیه ۲.۲. [۱۶] متغیرهای تصادفی X و Y با تکیه‌گاه‌های
 به ترتیب $\{x_1, \dots, x_N\}$ و $\{y_1, \dots, y_M\}$ و تابع g با مقادیر

$$f_\theta(x|t) = P_\theta(\mathbf{X} = \mathbf{x}|T = t)$$

$$= \begin{cases} \frac{1}{\binom{n}{t}}, & T(\mathbf{x}) = t \\ 0, & \text{در غیر این صورت} \end{cases} \quad (5)$$

بنا بر این $T(\mathbf{X}) \rightarrow (X_1, \dots, X_n) \rightarrow \theta$ از
 یک زنجیر مارکوف می‌آید و لذا $T(\mathbf{X})$ برای θ بسنده
 است. از دیدگاه دیگر نیز می‌توان بسندگی T را
 نتیجه گرفت. از آنجا که برای یک $\theta \in (0, 1)$ و
 $t \in \mathcal{T} = \{0, \dots, n\}$ دلخواه، همه x ها با t تا 1 تا $n-t$
 صفر دارای احتمال‌های برابر هستند و تعداد چنین x ها برابر با
 ترکیب $\binom{n}{t}$ است، رابطه (۵) در معادله (۴) صدق می‌کند که
 در آن $g(\mathbf{x}, t)$ توزیع یکنواخت روی همه x هایی است که دقیقاً
 t تا ۱ دارند و بنا بر این $T(\mathbf{X})$ یک آماره بسنده برای θ است.

یک آماره بسنده ممکن است حاوی اطلاعاتی باشد که
 مناسب نیست و برای رفع این مشکل از آماره بسنده مینیمال
 استفاده می‌شود که این آماره اطلاعات در داده‌ها راجع به پارامتر
 را به طور ماکسیمال فشرده می‌کند. آماره بسنده $T(X)$ را برای
 θ بسنده مینیمال گویند هرگاه تابعی از هر آماره بسنده دیگر مثل
 $U(X)$ باشد [۱۱]. بنا بر این $X \rightarrow U(X) \rightarrow T(X) \rightarrow \theta$ که
 در این حالت اطلاعات حاوی θ در نمونه بیشترین فشرده‌گی را
 دارد و همچنین

$$I(X, T(X)) \leq I(X, U(X)), \quad (6)$$

که این نامساوی نشان می‌دهد از بین همه توابع بسنده، بسندگی
 مینیمال دارای حداقل اطلاع متقابل روی نمونه X است. به کمک
 روابط (۳) و (۶) در حالتی که $(X, Y) \sim p(x, y)$ باشد، می‌توان
 قسمت وابسته X نسبت به Y را تعریف کرد. در این حالت
 متغیر تصادفی T را به گونه‌ای پیدا می‌کنیم که:

$$1. \quad Y \rightarrow X \rightarrow T \text{ از یک زنجیر مارکوف بیاید.}$$

۲. $I(X, T)$ کمینه باشد (مینیمال بودن)، در حالی
 که $I(Y, T)$ بیشینه باشد (بسندگی)، که این دو

^۶information bottleneck

با توجه به این حقیقت که $\mathcal{R} = \mathcal{U}$ معادل با این است که به ازای هر $j = 1, \dots, N$ و $k = 1, \dots, M$ داشته باشیم $r_{jk} = u_{jk}$ پس $\frac{r_{jk}}{p_j} = \frac{t_{v(x_j)k}}{t_{v(x_j)}}$ بنا بر این اگر $g(x_i) = g(x_j)$ باشد، از این که $v(x_i) = v(x_j)$ داریم:

$$\frac{r_{ik}}{p_i} = \frac{t_{v(x_i)k}}{t_{v(x_i)}} = \frac{t_{v(x_j)k}}{t_{v(x_j)}} = \frac{r_{jk}}{p_j},$$

یا به طور معادل

$$\frac{P(X = x_i, Y = y_k)}{P(X = x_i)} = \frac{P(X = x_j, Y = y_k)}{P(X = x_j)},$$

که همان رابطه (V) است.

از این رو ثابت کردیم رابطه (V) هم ارز با $\mathcal{R} = \mathcal{U}$ است. پس با استفاده از گزاره ۱.۱ و تساوی

$$I(X, Y) - I(g(X), Y) = D(\mathcal{R} \parallel \mathcal{U}),$$

شرط لازم و کافی برای $I(X, Y) = I(g(X), Y)$ آن است که $\mathcal{R} = \mathcal{U}$ باشد. بنا بر این رابطه (V) برقرار است اگر و تنها اگر $I(g(X), Y) = I(X, Y)$ باشد، که همان تعریف بسندگی $g(X)$ نسبت به متغیر تصادفی Y است. به علاوه اگر $g(x_j) = z_l$ آن گاه $\{X = x_j\} \subseteq \{g(X) = z_l\}$ و در پی آن

$$\begin{aligned} P(X = x_j, Y = y_k | g(X) = z_l) &= \frac{P(X = x_j, Y = y_k, g(X) = z_l)}{P(g(X) = z_l)} \\ &= \frac{P(X = x_j, Y = y_k)}{P(g(X) = z_l)} \\ &= \frac{P(X = x_j)P(Y = y_k | X = x_j)}{P(g(X) = z_l)} \\ &= \frac{P(X = x_j, g(X) = z_l)}{P(g(X) = z_l)} P(Y = y_k | X = x_j) \\ &= P(X = x_j | g(X) = z_l) P(Y = y_k | X = x_j), \end{aligned}$$

و لذا با مقایسه تساوی بالا با رابطه (A) می توان گفت به شرط $g(x_j) = z_l$ رابطه (A) معادل است با

$$P(Y = y_k | X = x_j) = P(Y = y_k | g(X) = z_l).$$

بنا بر این (V) و (A) معادل هستند. □

با توجه به تعریف ۲.۱، می توان یک تعریف معادل برای زنجیر مارکوف ارائه کرد.

حقیقی و بولر اندازه پذیر را در نظر بگیرید. شرط لازم و کافی برای آن که $g(X)$ یک تابع بسنده برای Y باشد آن است که توزیع احتمال شرطی Y به شرط X فقط به مقدار $g(X)$ بستگی داشته باشد؛ یعنی اگر به ازای $i = 1, \dots, N$ و $j = 1, \dots, N$ داشته باشیم $g(x_i) = g(x_j)$ آن گاه

$$P(Y = y_k | X = x_i) = P(Y = y_k | X = x_j). \quad (V)$$

به عبارت دیگر X و Y متغیرهای تصادفی مستقل هستند هنگامی که مقدار $g(X)$ ثابت باشد، یعنی برای هر z که $P(g(X) = z) > 0$ باشد، داریم:

$$\begin{aligned} P(X = x_j, Y = y_k | g(X) = z) &= \\ P(X = x_j | g(X) = z) P(Y = y_k | g(X) = z). \quad (A) \end{aligned}$$

اثبات. تعریف کنید

$$\mathcal{P} = \{p_1, \dots, p_N\}$$

$$p_j = P(X = x_j) \text{ که}$$

$$\mathcal{Q} = \{q_1, \dots, q_M\}$$

$$\text{که } q_k = P(Y = y_k)$$

$$\mathcal{R} = \{r_{jk}, \quad j = 1, \dots, N, \quad k = 1, \dots, M\}$$

که در آن

$$r_{jk} = P(X = x_j, Y = y_k).$$

فرض کنید $\{z_1, \dots, z_s\}$ مجموعه مقادیر مجزا از برد تابع g باشد که مقادیر خود را روی مجموعه $\{x_1, \dots, x_N\}$ اختیار می کنند و قرار دهید $g(x_j) = z_l$ که $j \in D_l$ ، D_1, \dots, D_s یک افراز از مجموعه $\{1, \dots, N\}$ است. به علاوه برای $j \in D_l$ تعریف کنید $l = v(x_j)$ و همچنین $t_{lk} = P(g(X) = z_l, Y = y_k)$ و $t_l = P(g(X) = z_l)$ اعداد

$$u_{jk} = \frac{t_{v(x_j)k} \cdot p_j}{t_{v(x_j)}}, \quad j = 1, \dots, N, \quad k = 1, \dots, M,$$

از یک توزیع احتمال $\mathcal{U} = \{u_{jk}\}$ می آیند؛ زیرا

$$\begin{aligned} \sum_{j=1}^N \sum_{k=1}^M u_{jk} &= \sum_{j=1}^N \sum_{k=1}^M \frac{t_{lk} \cdot p_j}{t_l} \\ &= \sum_{j=1}^N \sum_{k=1}^M \frac{P(g(X) = z_l, Y = y_k) P(X = x_j)}{P(g(X) = z_l)} \\ &= \sum_{j=1}^N P(X = x_j) \left[\frac{\sum_{k=1}^M P(g(X) = z_l, Y = y_k)}{P(g(X) = z_l)} \right] \\ &= 1. \end{aligned}$$

تعریف ۲.۳. تابع $T : \mathcal{X} \rightarrow \mathcal{T}$ را آمارهٔ بسندهٔ احتمالاتی برای $\{f_\theta\}$ گویند هرگاه برای هر $\theta \in \Theta$ داشته باشیم:

$$\sum_x f_\theta(x) \log\left(\frac{1}{f_\theta(x|T(x))}\right) = \sum_x f_\theta(x) \log\left(\frac{1}{g(x, T(x))}\right).$$

تعریف ۳.۳. تابع $T : \mathcal{X} \rightarrow \mathcal{T}$ را آمارهٔ بسندهٔ احتمالاتی برای $\{f_\theta\}$ گویند هرگاه برای هر توزیع پیشین $\pi(\theta)$ روی Θ داشته باشیم:

$$\sum_{\theta, x} \pi(\theta, x) \log\left(\frac{1}{f_\theta(x|T(x))}\right) = \sum_{\theta, x} \pi(\theta, x) \log\left(\frac{1}{g(x, T(x))}\right).$$

برای اثبات این حقیقت که تعریف ۱.۳، تعریف ۲.۳ را نتیجه می‌دهد فرض کنید به‌ازای هر $\theta \in \Theta$ و $x \in \mathcal{X}$ رابطه (۹) برقرار باشد. با گرفتن امید ریاضی از طرفین رابطه (۹)، تعریف ۲.۳ ثابت می‌شود. برای اثبات عکس، فرض کنید برای هر $\theta \in \Theta$ تعریف ۲.۳ برقرار باشد. با تعریف $f_\theta(t) = \sum_{y \in \mathcal{X}: T(y)=t} f_\theta(y)$ و اضافه کردن

$$\sum_x f_\theta(x) \log \frac{1}{f_\theta(T(x))},$$

به هر دو طرف تساوی، تعریف ۲.۳ را می‌توان دوباره به‌صورت زیر بازنویسی کرد:

$$\sum_x f_\theta(x) \log \frac{1}{f_\theta(x)} = \sum_x f_\theta(x) \log \frac{1}{g_\theta(x)}, \quad (10)$$

که $g_\theta(x) = f_\theta(T(x))g(x, T(x))$. با توجه به تعریف ۱.۱ رابطه (۱۰) زمانی برقرار است که برای هر $x \in \mathcal{X}$ داشته باشیم $g_\theta(x) = f_\theta(x)$ و در پی آن تعریف ۱.۳ نتیجه می‌شود.

برای اثبات معادل بودن تعریف ۲ و تعریف ۱.۳ از خطی بودن امید ریاضی و این که $p(\theta, x) = \pi(\theta)f_\theta(x)$ استفاده می‌شود. پیچیدگی کولموگوروف^۷ (اطلاع یا آنتروپی الگوریتمی) $K(x)$ ، برای اولین بار در مراجعی نظیر [۹، ۲، ۱۹] به‌عنوان طول کوتاه‌ترین توصیف دودویی مؤثر از x معرفی شد. به‌عبارت دیگر $K(x)$ ، طول کوتاه‌ترین برنامهٔ دودویی برای محاسبهٔ رشتهٔ x از یک کامپیوتر همگانی نظیر ماشین تورینگ^۸ است [۹]. ماشین تورینگ یک مدل ریاضی برای وسیله‌ای است که می‌تواند داده‌ها را بر یک نوار ذخیره‌سازی قابل کنترل بخواند و بنویسد. مثلاً $K(x)$ می‌تواند به‌عنوان طول کوتاه‌ترین برنامهٔ کامپیوتری باشد که x چاپ و سپس برنامهٔ متوقف شود. این برنامهٔ کامپیوتری می‌تواند با یکی از زبان‌های فراگیر کامپیوتری نظیر

تعریف ۳.۲. دنبالهٔ $\{X_n\}_{n \in \mathbb{N}}$ از متغیرهای تصادفی با مقادیر حقیقی روی فضای احتمال (Ω, \mathcal{F}, P) را که تعداد متناهی از مقادیر را اختیار می‌کند زنجیر مارکوف گسسته‌پارامتر گویند هرگاه برای $n \in \mathbb{N}$ داشته باشیم

$$I((X_1, \dots, X_n), X_{n+1}) = I(X_n, X_{n+1}).$$

نتیجهٔ زیر، ارتباط بین بسندگی و زنجیر مارکوف را بیان می‌کند.

نتیجه ۴.۲. تعریف ۲ بیان‌کنندهٔ این حقیقت است که برای یک زنجیر مارکوف گسسته‌پارامتر اطلاعی که راجع به X_{n+1} بر اساس مشاهدات X_1, \dots, X_n به دست می‌آید برابر با اطلاعی است که راجع به X_{n+1} بر اساس فقط مشاهدهٔ آخر، یعنی X_n ، به دست می‌آید و از این رو X_n که به‌صورت یک تابع از X_1, \dots, X_n است برای هر $n \in \mathbb{N}$ یک تابع بسنده برای X_{n+1} می‌باشد. به‌طور معادل،

$$P(X_{n+1} = x_{n+1} | X_1, \dots, X_n) = P(X_{n+1} = x_{n+1} | X_n).$$

۳ بسندگی احتمالاتی و بسندگی الگوریتمی

در ابتدا تعریفی برای آمارهٔ بسندهٔ احتمالاتی ارائه می‌کنیم. برای این منظور با نوشتن رابطه (۴) به‌صورت

$$\log \frac{1}{g(x, t)} = \log \frac{1}{f_\theta(x|t)}, \quad (9)$$

می‌توان تعریف زیر را ارائه کرد.

تعریف ۱.۳. تابع $T : \mathcal{X} \rightarrow \mathcal{T}$ را آمارهٔ بسندهٔ احتمالاتی برای $\{f_\theta\}$ گویند هرگاه تابع $g : \mathcal{X} \times \mathcal{T} \rightarrow \mathbb{R}$ وجود داشته باشد به‌طوری که برای هر $\theta \in \Theta$ و $x \in \mathcal{X}$ و $t \in \mathcal{T}$ در شرط (۹) صدق کند. (برای راحتی قرارداد می‌کنیم $\log \frac{1}{\cdot} = \infty$)

برای به دست آوردن نسخهٔ امید ریاضی تعریف بسندگی احتمالاتی، ثابت می‌شود تعریف‌های زیر هم‌ارز با تعریف ۱.۳ هستند.

^۷ Kolmogorov complexity

^۸ Turing machine

(۳) به ازای هر m مجموعه $\{s : \exists x \in \{0, 1\}^n : S(x) = s\}$ یک افراز از $\{0, 1\}^n$ باشد.

دنباله‌ای از توابع جرم احتمال $\{f_\theta^{(n)}\}_{n \in \mathbb{N}}$ را که $f_\theta^{(n)}$ تابع جرم احتمال روی $\{0, 1\}^n$ است و توزیع کناری f_θ روی اولین قطعه n بیتی است، در نظر بگیرید. آماره دنباله‌ای S را تقریباً بسنده برای θ از دیدگاه احتمالاتی-فردی گویند هرگاه توابع $g^{(1)}$ و $g^{(2)}$ و ثابت c وجود داشته باشند به قسمی که به ازای هر θ و n و $x \in \{0, 1\}^n$ داشته باشیم:

$$\left| \log \frac{1}{f_\theta^{(n)}(x|S(x))} - \log \frac{1}{g^{(n)}(x|S(x))} \right| < \varepsilon,$$

و S را تقریباً بسنده برای θ از دیدگاه احتمالاتی-امید ریاضی گویند هرگاه توابع $g^{(1)}$ ، $g^{(2)}$ ، ... و ثابت c' وجود داشته باشند به طوری که به ازای هر θ و n داشته باشیم:

$$\left| \sum_{x \in \{0, 1\}^n} f_\theta^{(n)}(x) \left[\log \frac{1}{f_\theta^{(n)}(x|S(x))} - \log \frac{1}{g^{(n)}(x|S(x))} \right] \right| \leq c'. \quad (11)$$

آماره دنباله‌ای S را از دیدگاه الگوریتمی، بسنده گویند هرگاه ثابت c وجود داشته باشد که برای هر n و $x \in \{0, 1\}^n$ برنامه مولد $S(x)$ یک آماره بسنده الگوریتمی برای x (نسبت به ثابت c) باشد؛ یعنی

$$K(S(x)) + \log|S(x)| \leq K(x) + c.$$

در قضیه زیر نشان می‌دهیم که آماره بسنده الگوریتمی، آماره بسنده احتمالاتی نیز هست.

قضیه ۶.۳. اگر S یک آماره دنباله‌ای باشد که از لحاظ الگوریتمی بسنده است، به ازای هر θ که $K(f_\theta) < \infty$ ، یک ثابت c وجود دارد که برای هر n ، نامساوی (۱۱) برقرار است و در آن توزیع یکنواخت است؛ یعنی

$$g^{(n)}(x|s) = \begin{cases} \frac{1}{|\{x \in \{0, 1\}^n : S(x) = s\}|}, & S(x) = s \\ 0, & \text{در غیر این صورت} \end{cases}$$

C ، جاوا^۱ و $LISP$ ^{۱'} نوشته شود و در هر یک از آنها نتیجه برنامه برای مقدار $K(x)$ در یک ثابت که به بستگی ندارد متفاوت است. به طور مشابه پیچیدگی کولموگوروف شرطی $K(x|y)$ طول کوتاه‌ترین برنامه برای محاسبه x است، به شرطی که y یک ورودی جانبی (مثل نوار مغناطیسی، فلاپی دیسک، فلش مموری و لوح فشرده) باشد [۶]. برای بررسی دقیق‌تر موضوع، فرض کنید T_1, T_2, \dots شمارش استاندارد همه ماشین‌های تورینگ باشد و ϕ_1, ϕ_2, \dots شمارش توابع متناظر که توسط یک ماشین تورینگ مخصوص محاسبه شده است. بنا بر این T_i مقدار ϕ_i را محاسبه می‌کند. مجموعه S شامل x با ویژگی $K(x) = K(S) + \log|S| + O(1)$ را بهینه نامند. یک آماره بسنده الگوریتمی x کوتاه‌ترین برنامه برای یک مجموعه S حاوی x است که بهینه باشد، یعنی در شرط معادله بالا صدق کند. یک آماره بسنده الگوریتمی با مجموعه بهینه S را مینیمال گویند هرگاه هیچ مجموعه بهینه‌ای مثل S' وجود نداشته باشد که $K(S') < K(S)$. برای بهتر درک کردن مفهوم مجموعه بهینه، مثال زیر می‌تواند مفید باشد.

مثال ۴.۳. فرض کنید $k \in \{0, 1, \dots, n\}$ باشد و x یک رشته به طول n با تعداد k تا ۱ باشد. این x می‌تواند به عنوان نتیجه پرتاب یک سکه با اربیبی $p = \frac{k}{n}$ در نظر گرفته شود. توصیف دو قسمتی از x را می‌توان به صورت تعداد ۱ها در x برای اولین بار (k) که اندیس در شرط $j \leq \log|S|$ صدق می‌کند و x در یک مجموعه S حاوی رشته‌های به طول n با k تا ۱ است، در نظر گرفت. این مجموعه بهینه است، زیرا $K(x|n) = K(S) + \log|S|$.

اکنون به ارتباط بسندگی الگوریتمی و بسندگی احتمالاتی می‌پردازیم. برای این منظور، تعریف زیر ارائه می‌شود.

تعریف ۵.۳. یک آماره دنباله‌ای، تابعی مانند $S : \{0, 1\}^* \rightarrow \{0, 1\}^*$ است (که در آن

$$\{0, 1\}^* = \{\varepsilon, 0, 1, 00, 01, 10, 11, 000, \dots\}$$

و ε لغت تهی (حرفی نداشته باشد) است) به طوری که برای هر n و $x \in \{0, 1\}^n$ داریم:

$$S(x) \subseteq \{0, 1\}^n \quad (1)$$

$$x \in S(x) \quad (2)$$

^۱ Java

^{۱'} locator identifier separation protocol

با استفاده از گزاره ۱.۱ داریم:

$$\sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) K(S(x)) \geq \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) \log \frac{1}{f_\theta^{(n)}(S(x))}.$$

توجه کنید که برای هر n داریم:

$$\begin{aligned} & \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) \left[\log \frac{1}{f_\theta^{(n)}(S(x))} + \log \frac{1}{f_\theta^{(n)}(x|S(x))} \right] \\ &= \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) \log \frac{1}{f_\theta(x)}, \end{aligned} \quad (14)$$

و دو بخشی را که x توسط اولین رمزگذاری $S(x)$ با استفاده از $\log \frac{1}{f_\theta^{(n)}(S(x))}$ بیت و سپس توسط $\log |S(x)|$ بیت رمزنگاری می‌شود در نظر بگیرید. با توجه به گزاره ۱.۱ این کد باید کارایی کمتری نسبت به کد شانون-فانو با طول $\log \frac{1}{f_\theta(x)}$ داشته باشد. بنا بر این با توجه به (۱۴) برای هر n داریم:

$$\sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) \log |S(x)| \geq \log \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) \log \frac{1}{f_\theta^{(n)}(x|S(x))}.$$

بنا بر این یک c' وجود دارد که برای هر n داشته باشیم:

$$\left| \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) \log |S(x)| - \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) \log \frac{1}{f_\theta^{(n)}(x|S(x))} \right| \leq c'.$$

□

۴ نتیجه گیری

در این مقاله با تعریف مفهوم بسندگی یک آماره نسبت به یک متغیر تصادفی دیگر از دیدگاه نظریه اطلاع، یک ارتباط بین مفهوم بسندگی و زنجیره‌های مارکوف نتیجه می‌شود. ذکر این نکته جالب است که در حالت خاص که یک آماره نسبت به یک متغیر تصادفی ثابت (که در این جا مقدار پارامتری را اختیار می‌کند که باید برآورد شود) بسنده باشد همان مفهوم بسندگی معمولی نتیجه می‌شود. علاقه‌مندان با در نظر گرفتن پارامتر به‌عنوان یک متغیر تصادفی می‌توانند این مفهوم را در حالت بیزی گسترش دهند، که می‌تواند شروعی برای تحقیقات آینده باشد. در بخش آخر مقاله مفهوم بسندگی الگوریتمی مطرح می‌شود و ثابت می‌شود که همان مفهوم بسندگی احتمالاتی را نتیجه می‌دهد.

بنا بر این اگر $\sup_{\theta \in \Theta} K(f_\theta) < \infty$ باشد، S یک آماره تقریباً بسنده برای θ از دیدگاه احتمالاتی-امید ریاضی است، که g توزیع یکنواخت است.

اثبات. طبق تعریف بسندگی الگوریتمی، یک ثابت c وجود دارد به‌طوری که به‌ازای هر θ و n داریم:

$$\begin{aligned} & \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) [K(S(x)) + \log |S(x)|] \\ & \leq \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) K(x) + c. \end{aligned} \quad (12)$$

اکنون با ثابت فرض کردن θ که $K(f_\theta) < \infty$ ، به‌ازای یک $c_\theta \approx K(f_\theta)$ و هر n خواهیم داشت:

$$\cdot \leq \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) K(x) - \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) \log \frac{1}{f_\theta(x)} \leq c_\theta, \quad (13)$$

که دلیل نامساوی اولی، گزاره ۱.۱ است. دلیل نامساوی دوم این است که با توجه به $K(f_\theta) < \infty$ ، کد شانون-فانو^{۱۱} می‌تواند توسط یک برنامه کامپیوتری با یک اندازه ثابت انجام شود که به n بستگی ندارد.

با توجه به (۱۲) و (۱۳) برای هر n داریم:

$$\begin{aligned} & \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) \log \frac{1}{f_\theta(x)} \\ & \leq \sum_{x \in \{0,1\}^n} f_\theta^{(n)}(x) [K(S(x)) + \log |S(x)|] \\ & \leq \sum_x f_\theta^{(n)}(x) \log \frac{1}{f_\theta(x)} + c_\theta. \end{aligned}$$

برای $s \subseteq \{0,1\}$ نماد

$$f_\theta^{(n)}(s) = \sum_{y \in \mathcal{X}: S(y)=s} f_\theta^{(n)}(y)$$

را در نظر می‌گیریم. با توجه به فرضیات قسمت ۳ از تعریف ۵.۳ داریم:

$$\sum_{s: \exists x \in \{0,1\}^n: S(x)=s} f_\theta^{(n)}(s) = 1,$$

که $f_\theta^{(n)}(s)$ تابع جرم احتمال روی S شامل مجموعه مقادیر آماره S است که می‌تواند دنباله‌ای به طول n را اختیار کند. بنا بر این

^{۱۱} Shannon-Fano code

مراجع

- [1] Basharin, G. P., Langville, A. N. and Naumov, V. A. (2004). The life and work of A. A. Markov. *Linear Algebra and its Applications*, **386**, 3-26.
- [2] Chaitin, G. J. (1969). On the length of programs for computing finite binary sequences: statistical considerations. *journal of Association for Computing Machinery*, **16**, 145-159.
- [3] Cover, T. M. and Thomas, J. A. (2006). *Elements of Information Theory*. 2nd edition, Wiley.
- [4] Fushing, H., Wang, H., Vander Waal, K., McCowan, B. and Koehl, P. (2013). Multi-scale clustering by building a robust and self correcting. *Ultrametric Topology on Data Points, PloS one*, **8**.
- [5] Fushing, H., Kevin, F. and Cho-Jui, H. (2016). Machine learning meliorates computing and robustness in discrete combinatorial optimization problems. *Frontiers in Applied Mathematics and Statistics*, **2**, 1-8.
- [6] Gacs, P., Tromp, J. and Vitanyi, P. M. B. (2001). Algorithmic statistics. *IEEE Transactions on Information Theory*, **47**, 2443-2463.
- [7] Ghurye, S. G. (1968). Information and sufficient sub-fields. *Annals of Mathematical Statistics*, **39**, 2056-2066.
- [8] Gray, R. M. (2009). *Entropy and information theory*. Springer.
- [9] Kolmogorov, A. N. (1965). Three approaches to the quantitative definition of information. *Problems of Information Transmission*, **1**, 1-7.
- [10] Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *Annals of Mathematical Statistics*, **22**, 79-86.
- [11] Lehmann, E. L. and Romano, J. P. (2005). *Testing Statistical Hypotheses*. 3rd edition, Springer, New York.
- [12] Li, M. and Vitanyi, P. M. B. (1997). *An introduction to Kolmogorov complexity and its applications*. 2nd ed. New York: Springer-Verlag.
- [13] Longest, K. C. and Vaisey, S. (2008). fuzzy: A program for performing qualitative comparative analyses (QCA) in stata. *The Stata Journal*, **8**, 79-104.
- [14] Markov, A. A. (1907). Issledovanie zamechatel'nogo sluchaya zavisimyh ispytaniy. *Izvestiya Akademii Nauk, SPb, VI seriya*, **1(93)**, 61-80.
- [15] Markoff, A. (1910). Recherches sur un cas remarquable d'épreuves dépendantes. *Acta mathematica*, **33(1)**, 87-104.
- [16] Pop-Stojanovic, Z. R. (2006). A classroom note: entropy, information, and Markov property. *Teaching Mathematics*, **9**, 1-12.

- [17] Ramu, P., Qu, X., Youn, B., Haftka R. T., and Choi, K. (2006). Safety factors and inverse reliability measures. *International Journal of Reliability and Safety*, **1**, 1-2.
- [18] Shannon, C. E. (1948). The mathematical theory of communication. *The Bell System Technical Journal*, **27**, 623-656.
- [19] Solomonoff, R. J. (1964). A formal theory of inductive inference, part 1 and part 2. *Information and computation*, **7**, 1–22.
- [20] Tishby, N., Pereira, F. C. and Bialek, W. (1999). The information Bottleneck method. in proceeding of the 37th Annual Allerton Conference on Communication, *Control and Computing*, 368-377.
- [21] Vereshchagin, N. K. (2009). *Algorithmic minimal sufficient statistic Revisited*. In CiE (pp. 478-487).
- [22] Vereshchagin, N. (2016). Algorithmic minimal sufficient statistics: a new approach. *Theory of Computing Systems*, **58**, 463-481.
- [23] Von Hilgers, P. and Langville, A. N. (2006). *The Five Greatest Applications of Markov Chains*. Proceedings of the Markov Anniversary Meeting. Bosen Press.
- [24] Wallace, C. and Freeman, P. (1987). Estimation and inference by compact coding. *Journal of the Royal Statistical Society*, **49**, 240–251.