

مجله علوم آماری، پاییز و زمستان ۱۳۹۴

جلد ۹، شماره ۲، ص ۲۵۳-۲۶۷

مدل‌های بیزی در انطباق پروتئین‌ها

سید مرتضی نجیبی^۱، موسی گل‌علی‌زاده^۲، محمدرضا فقیهی^۱

^۱ گروه آمار، دانشگاه شهید بهشتی

^۲ گروه آمار، دانشگاه تربیت مدرس

تاریخ دریافت: ۱۳۹۳/۱۰/۲ تاریخ آخرین بازنگری: ۱۳۹۴/۲/۲۸

چکیده: در این مقاله راه‌حل‌های احتمالاتی مسئله انطباق ساختار سوم پروتئین و تفاوت آن با الگوریتم‌های قطعی مطالعه و بررسی می‌شود. برای این منظور دو مدل احتمال بیزی معرفی و راه‌حلی در خصوص اضافه کردن اطلاعات توالی و نوع اسید آمینه (ساختار اول) به آن ارائه خواهد شد. همچنین نحوه برآورد پارامترهای انطباق به کمک الگوریتم مونت کارلوی زنجیر مارکوف و نمونه‌گیری از توزیع پسین معرفی می‌شود. در نهایت نحوه کاربست این روش‌ها در انطباق پروتئین‌ها نشان داده شد و برآورد پارامترها در مدل‌های متفاوت برای یک مجموعه داده واقعی ارزیابی و مقایسه خواهد شد.

واژه‌های کلیدی: آمار شکل، فاصله اندازه-شکل، بیوانفورماتیک، انطباق ساختاری پروتئین، ساختار سوم پروتئین، استنباط بیزی.

آدرس الکترونیک مسئول مقاله: سید مرتضی نجیبی، sm_najibi@sbu.ac.ir
کد موضوع‌بندی ریاضی (۲۰۱۰): ۶۲H۱۱، ۶۲P۱۰، ۶۵D۱۸

تاکنون الگوریتم‌های بسیاری در حوزه بیوانفورماتیک ساختاری به منظور رده‌بندی ساختار هندسی پروتئین‌ها معرفی شده‌اند که عموماً بر پایه انطباق زوجی شکل هندسی پروتئین‌ها هستند. هدف اصلی این روش‌ها، انطباق پروتئین‌ها بر یکدیگر به گونه‌ای است که بعد از دوران و تغییر مکان، فاصله نقاط جور شده از یکدیگر حداقل و تعداد نقاط جور شده ماکسیمم باشد. از آنجا که بررسی تمامی حالت‌های ممکن امکان‌پذیر نیست، تقریباً در تمامی روش‌ها از روش‌های اکتشافی^۱ برای دستیابی به مدل بهینه استفاده شده است. اما در عمل این الگوریتم‌ها لزوماً جواب یکسانی برای یک مسئله ارائه نمی‌دهند. به عنوان مثال میزان شباهت دو روش مرسوم انطباق ماتریسی^۲ (DALI) و تعمیم ترکیبی^۳ (CE) حداکثر ۴۰ درصد است (شیندیالو و بورنه، ۲۰۰۰). از این رو به نظر می‌رسد مقایسه آماری پروتئین‌ها تنها به کمک فاصله آن‌ها از یکدیگر، اطلاعات واقعی راجع به ساختار هندسی شان را حذف کرده و در نتیجه جواب نهایی تحلیل، رضایت‌بخش نیست. لذا انتظار می‌رود در نظر گرفتن ساختار هندسی پروتئین به عنوان منبع اطلاعات ذی‌قیمت آماری بتواند در ارزیابی میزان دقیق مشابهت و تفاوت ساختمان آن‌ها مفید و موثر واقع شود.

برای مقایسه اشکال هندسی در آمار شکل، مدل‌های آماری متنوعی معرفی شده‌اند (درایدن و ماردیا، ۱۹۹۸). این روش‌ها در مقایسه با روش‌های قطعی، از رویکرد احتمالاتی استفاده می‌کنند و در یک دید کلی هدف‌شان مدل‌سازی خطاهای تصادفی است (تیلور و همکاران، ۲۰۰۳). اطلاعات اولیه آماری شکل در این الگوریتم‌ها به کمک نقاط شاخص^۴ به دست می‌آید. اما در مسائلی همچون انطباق پروتئین‌ها، تناظر بین نقاط یا برچسب^۵ آن‌ها مشخص نیست. از این رو برخی از محققان، مانند کنت و همکاران (۲۰۰۴)، گرین و ماردیا (۲۰۰۶)، درایدن و همکاران

۱ Heuristic
 ۲ Distance alignment matrix
 ۳ Combinatorial extension
 ۴ Landmarks
 ۵ Label

(۲۰۰۷)، کنوبی و درایدن (۲۰۱۲) و ماردیا و همکاران (۲۰۱۳) مطالعاتی را در زمینه تحلیل اشکال بدون برجسب انجام داده‌اند. اما استفاده از این روش‌ها در انطباق ساختاری پروتئین، به راحتی و سرعت الگوریتم‌های قطعی نیست. یکی از راه‌حل‌های بهبود این مدل‌ها، اضافه کردن اطلاعات اضافی پروتئین‌ها در قالب توزیع پیشین آگاهی‌بخش است (فالیزی و همکاران، ۲۰۱۴؛ نجیبی و همکاران، ۲۰۱۵).

در این مقاله دو مدل احتمالاتی انطباق زوجی اشکال بدون برجسب معرفی و نحوه اضافه کردن اطلاعات مربوط به توالی اسیدهای آمینه و نوع آن‌ها به کمک پیشین‌های آگاهی‌بخش تجربی ارائه خواهد شد. همچنین کاربرد این مدل‌ها در یک مجموعه داده واقعی و میزان تأثیر اضافه کردن اطلاعات ساختاری پروتئین به مدل‌های قبلی مورد مطالعه قرار گرفته است.

۲ مدل‌های بیزی انطباق پروتئین‌ها

فرض کنید مجموعه X با k نقطه و مجموعه μ با l نقطه ($k \neq l$) در فضای m بعدی ($m \geq 2$) در اختیار باشند. در اصطلاح آمار شکل به این دو مجموعه پیکره یا پیکربندی و به فضایی که مشاهدات در آن قرار می‌گیرند، فضای پیکربندی گویند (درایدن و ماردیا، ۱۹۹۸). از آنجا که برجسب گذاری بین نقاط نامعلوم است، مسئله انطباق پروتئین در واقع یافتن زیرمجموعه‌هایی از دو پیکره X و μ تعریف می‌شود به گونه‌ای که بعد از انجام تبدیلات انتقال^۶ و دوران، فاصله این دو زیرمجموعه از یکدیگر حداقل شود. این فاصله در واقع اندازه‌ای است که در فضای مشاهدات بعد از انجام تبدیلات انتقال و دوران تعریف می‌شود.

در این مسئله مختصات اتم‌های کربن-آلفا در هر پروتئین، نقاط شاخص در دو پیکره X و μ هستند. تناظر بین این نقاط از طریق برآورد ماتریسی به نام ماتریس انطباق^۷ (A) به دست می‌آید که دارای k سطر و $l + 1$ ستون است. حداکثر یکی از عناصر هر سطر این ماتریس یک و مابقی صفر هستند. ستون $l + 1$ این ماتریس

^۶ Translation

^۷ Match matrix

نشان دهنده نقاطی از X است که با هیچ نقطه‌ای از μ جور نشده‌اند. بنابراین بر اساس این ماتریس، در صورتی که فرض شود تعداد p نقطه از دو پیکربندی با یکدیگر جور شده‌اند، می‌توان ماتریس‌های X و μ را به صورت

$$X = (X^M, X^N), \quad \mu = (\mu^M, \mu^N)$$

افراز کرد، که در آن $X_{p \times m}^M$ و $\mu_{p \times m}^M$ به ترتیب دو زیرمجموعه از نقاط X و μ هستند که با یکدیگر جور شده‌اند و $X_{(k-p) \times m}^N$ و $\mu_{(\ell-p) \times m}^N$ نیز دو زیرمجموعه از نقاط X و μ هستند که با یکدیگر جور نشده‌اند.

مدل‌های انطباق بیزی این مسئله بر اساس نحوه برآورد تبدیلات انتقال و دوران، به دو گروه کلی تقسیم‌بندی می‌شوند. در گروه اول، مطابق لی (۱۹۸۸) و درآیدن و همکاران (۲۰۰۷)، مدل انطباق در فضای غیر اقلیدسی اندازه-شکل تعریف می‌شود که در آن بردار انتقال γ و ماتریس دوران Γ به کمک تحلیل پروکراستس به دست می‌آیند و جزء پارامترهای اصلی مدل محسوب نمی‌شوند. در گروه دوم، مدل‌بندی احتمالاتی انطباق بر اساس ماردیا (۲۰۰۶) و کنوبی و همکاران (۲۰۱۲) به‌طور مستقیم بر روی نقاط شاخص در فضای پیکربندی تعریف می‌شود و بردار انتقال γ و ماتریس دوران Γ جزء پارامترهای اصلی مدل هستند. در ادامه این مدل‌ها معرفی و در بخش سوم نحوه بهبود آنها مورد مطالعه قرار گرفته است.

۱.۲ مدل احتمال در فضای اندازه-شکل

بدون از دست دادن کلیت مسئله فرض کنید پیکربندی $X_{k \times 3}$ تصادفی و پیکربندی $\mu_{l \times 3}$ ثابت است به نحوی که p نقطه شاخص از پیکربندی X با p نقطه از μ ، نظیر به نظیر جور شده‌اند. به‌علاوه فرض کنید دو پیکربندی $X = C_p X$ و $\mu = C_p \mu$ ، قبلاً توسط ماتریس مرکزی کننده $C_p = I_p - \mathbf{1}_p \mathbf{1}_p^t / p$ مرکزی شده‌اند، که در آن I_p ماتریس همبندی p بعدی و $\mathbf{1}_p$ برداری شامل p مقدار یک است. توجه شود که $\mathbf{1}_p^t X = \mathbf{0}_m$ و $\mathbf{1}_p^t \mu = \mathbf{0}_m$ ، که در آن $\mathbf{0}_m$ برداری شامل p مقدار صفر است. مجموعه p نقطه در فضای سه بعدی که تحت تبدیلات انتقال و دوران پایا است،

مانند یک نقطه در فضای اندازه-شکل است که با $S\Sigma_p^p$ نمایش داده می شود. اگر اندازه مرکزی ماتریس X به صورت $S(X) = \|X\| = \sqrt{X^t X}$ در نظر گرفته شود، فاصله اندازه-شکل که در واقع متر ریمانی بین دو نقطه X و μ در فضای $S\Sigma_p^p$ است (لی، ۱۹۸۸)، به صورت

$$d_S^2(X^M, \mu^M) = \inf_{\Gamma \in SO(m), \gamma \in R} \|(\mu^M)^t - \Gamma(X^M)^t - \gamma \times \mathbf{1}_p\|^2$$

تعریف می شود، که در آن ماتریس دوران $\Gamma_{m \times m}$ دوران $(|\Gamma| = 1$ و $\Gamma \Gamma^t = \Gamma^t \Gamma = I_p)$ بردار انتقال و $SO(m)$ فضای ماتریس های دوران (تبدیلات متعامد مخصوص) هستند (کندل و همکاران، ۱۹۹۹). جواب کمترین توان های دوم این فاصله، همان فاصله پروکراستس خواهد شد. اما از آنجا که در پروتئین ها نقاط شاخص در فضای سه بعدی قرار دارند، به منظور برآورد این فاصله در فضای اندازه-شکل از تقریب آن در فضای مماس استفاده می شود (کنت و ماردیا، ۲۰۰۱). بنابراین می توان از توزیع نرمال چند متغیره در فضای مماس بر اندازه-شکل μ ، برای پرتیابی نقاط جور شده به صورت (درآیدن و همکاران، ۲۰۰۷)

$$f(X^M) \propto (2\pi)^{-\frac{Q}{2}} \tau^{\frac{Q}{2}} \exp\left\{-\frac{\tau}{2} d_S^2(X^M, \mu^M)\right\}$$

استفاده نمود، که در آن $d_S^2(X^M, \mu^M)$ برآورد فاصله اندازه-شکل بین نقاط جور شده در دو پیکربندی در فضای مماس، $\tau = 1/\sigma^2$ پارامتر دقت و $Q = pm - m(m-1)/2 - m$ بعد فضای اندازه-شکل است. در مقابل فرض می شود که $k-p$ نقطه جور نشده X^N به صورت مستقل دارای تابع چگالی یکنواخت در فضای کراندار A با حجم $|A|$ هستند، یعنی

$$f(X^N) \propto |A|^{k-p}, \quad X^N \in A^{k-p}.$$

در این صورت تابع درستنمایی پیکربندی X به شرط معلوم بودن A ، μ و τ به صورت

$$\begin{aligned} L(X|A, \mu, \tau) &\propto f(X^M|A, \tau, \mu) f(X^N|A) \\ &\propto |A|^{k-p} (2\pi)^{-Q/2} \tau^{Q/2} \exp\left\{-\frac{\tau}{2} d_S^2(X^M, \mu^M)\right\} \end{aligned}$$

خواهد بود. دو پارامتر τ و Λ پارامترهای اصلی این مدل‌بندی هستند. فرض کنید توزیع پیشین پارامتر τ به صورت $\Gamma(\alpha_0, \beta_0)$ و احتمال پیشین برای هر سطر ماتریس Λ به صورت

$$P(\lambda_{i(\ell+1)} = 1) = \psi \quad P(\lambda_{ij} = 1) = \frac{1 - \psi}{\ell} \quad 1 \leq j \leq \ell, \quad i = 1, \dots, k$$

تعریف شود، که در آن $0 \leq \psi \leq 1$. بنابراین چگالی پسین توأم τ و Λ به شرط مشاهدات \mathbf{X} و μ به صورت

$$\pi(\tau, \Lambda | \mathbf{X}, \mu) = \frac{\pi(\tau)\pi(\Lambda)L(\mathbf{X}|\Lambda, \mu, \tau)}{\sum_{\Lambda} \int_0^1 \pi(\tau)\pi(\Lambda)L(\mathbf{X}|\Lambda, \mu, \tau)}$$

خواهد شد. به دلیل عدم وجود صورت بسته برای توزیع پسین و همچنین به علت ساختار خاص پارامترها، امکان استنباط مستقیم بر اساس توزیع پسین وجود ندارد. لذا درآیدن و همکاران (۲۰۰۷) پیشنهاد دادند که با استفاده از روش‌های شبیه‌سازی $MCMC$ و الگوریتم‌های گیبز و متروپولیس-هستینگز، برآورد ماکسیمم توزیع پسین^۸ (MAP) پارامترها انجام شود. برای این منظور نیاز است توزیع‌های شرطی کامل دو پارامتر τ و Λ تعیین شود. به سادگی می‌توان دریافت که توزیع شرطی کامل τ به شرط بقیه پارامترها به صورت

$$\tau | (\mathbf{X}, \Lambda, \mu) \sim \Gamma\left(\alpha_0 + \frac{Q}{\tau}, \beta_0 + \frac{1}{\tau} d_S^2(\mathbf{X}^M, \mu^M)\right).$$

خواهد بود. بنابراین می‌توان به راحتی با تولید نمونه از توزیع پسین شرطی کامل τ در هر مرحله از شبیه‌سازی با قدم‌های گیبز این پارامتر را به‌روز کرد. به علاوه توزیع شرطی کامل Λ به صورت

$$\begin{aligned} \pi(\Lambda | \mathbf{X}, \mu, \tau) &\propto \pi(\Lambda)L(\mathbf{X}|\Lambda, \mu, \tau) \\ &\propto |\Lambda|^{k-p} (\tau\pi)^{-Q/2} \tau^{Q/2} \exp\left\{-\frac{\tau}{\tau} d_S^2(\mathbf{X}^M, \mu^M)\right\} \end{aligned}$$

است. اما به دلیل عدم وجود صورت بسته برای توزیع این پارامتر، شبیه‌سازی از آن به کمک الگوریتم متروپولیس-هستینگز با توزیع نامزد مناسب انجام می‌شود. درآیدن

^۸ Maximum A Posteriori

و همکاران (۲۰۰۷) پیشنهاد دادند که در هر مرحله از الگوریتم، سطر i ام از بین k سطر ماتریس انطباق به تصادف انتخاب و سپس ستون j از بین $l + 1$ ستون ماتریس A با احتمال‌های $(\frac{1-\psi}{l}, \dots, \frac{1-\psi}{l}, \psi)$ انتخاب و مقدار λ_{ij} برابر یک در نظر گرفته شود. سپس ماتریس انطباق جدید با احتمال

$$\alpha_A = \min\left\{1, \frac{\pi(A^*|X, \mu, \tau)q}{\pi(A|X, \mu, \tau)q^*}\right\}$$

پذیرفته شود، که در آن

$$\frac{q}{q^*} = \begin{cases} \frac{\psi}{\sqrt{l}} & \text{اگر نقطه جورنشده به نقطه جورشده تبدیل شود} \\ \frac{\sqrt{l}}{\psi} & \text{اگر نقطه جورشده به نقطه جورنشده تبدیل شود} \\ 1 & \text{اگر نقطه جور شده با نقطه دیگری در } \mu \text{ جور شود} \end{cases}$$

بنابراین بر اساس چنین توزیع پیشنهادی می‌توان نمونه‌های متفاوتی را از چگالی شرطی کامل A تولید نمود. چون در ساختار این الگوریتم، به روزرسانی پارامترها با استفاده از قدم‌های گیبز و متروپولیس هستینگز بوده و هر دو مرحله نیازمند رسیدن به همگرایی مناسب هستند، ممکن است تعداد تکرارهای زیادی برای همگرایی نیاز باشد. به علاوه تحلیل پروکراستس جزئی در فضای مماس در هر مرحله از الگوریتم به منظور برآورد پارامترهای مکان و دوران انجام می‌شود. بنابراین به طور طبیعی حجم محاسبات در هر مرحله زیاد و سرعت این الگوریتم کم است.

برای افزایش سرعت همگرایی این الگوریتم، کنویبی و درایدن (۲۰۱۲) پیشنهاد دادند که از چهار گام اضافی نزدیکی، دوران، انتقال و تلنگر^۹ (پرش‌های بزرگ) به منظور دستیابی به نقطه شروع مناسب، استفاده شود. الگوریتم انطباق در این حالت مشابه الگوریتم قبلی است با این تفاوت که در تکرارهای اولیه، پس از تکرار مشخصی (مثلاً N_{settle}) به تصادف یکی از پرش‌های بزرگ یک بار تکرار می‌شود. با انجام هر پرش بزرگ ماتریس انطباق با تغییر چندگانه نقاط جور شده به ماتریس جدیدی تبدیل خواهد شد و با تکرار الگوریتم به صورت معمول به تعداد N_{settle} تأثیر ناگهانی این پرش بزرگ به تعادل می‌رسد. این روند تا تکرار مشخصی از الگوریتم (مثلاً N_{init}) ادامه می‌یابد. پس از آن روند طبیعی الگوریتم ادامه پیدا

^۹ Flip

می‌کند تا به همگرایی مطلوب دست یابد. این تکنیک در مطالعات شبیه‌سازی و مثال‌های کاربردی کارایی بهتری نسبت به الگوریتم قبلی دارد. اما همچنان به دلیل وابستگی شدید به نقطه شروع، از کارایی لازم در پروتئین‌ها برخوردار نیست. لذا در بخش‌های بعدی راه‌حل‌های مناسبی برای حل این مشکل پیشنهاد می‌شود.

۲.۲ مدل احتمال در فضای پیکربندی

گرین و ماردیا (۲۰۰۶) یک مدل آماری مشابه در فضای پیکربندی بجای فضای اندازه-شکل، برای مسئله انطباق اشکال بدون برچسب ارائه داده‌اند. در مدل آن‌ها مسئله انطباق به یافتن جواب بهینه در رابطه

$$(\mu^M)^t = \Gamma(\mathbf{X}^M)^t + \gamma \times \mathbf{1}_p + \Sigma$$

تبدیل شد، که در آن ماتریس خطای $m \times p$ است و دو عنصر انتقال و دوران نیز به‌عنوان پارامترهای مدل محسوب شده و مانند بقیه پارامترها برآورد می‌شوند. در این روش، مدل‌سازی نقاط جور شده و جور نشده با استفاده از فرآیند پواسن و یک مدل بیزی به‌طور مستقیم بر روی داده‌ها انجام می‌شود. به‌طور دقیق‌تر، فرض می‌شود که هر دو پیکربندی $\mathbf{X}_k \times 3$ و $\mu_{\ell} \times 3$ ، از یک فرآیند پواسن با میانگین متغیر پنهان تبعیت می‌کنند و در نهایت با تشکیل تابع درستی‌نمایی برآورد پارامترهای این مدل به همراه ماتریس انطباق انجام می‌شود. با بهره‌گیری از این روش، کنوبی و درایدن (۲۰۱۲) فرض کردند که نقاط جور شده \mathbf{X}^M به‌صورت یک پرشیدگی تصادفی نرمال از میانگین μ^M باشند و نقاط جور نشده \mathbf{X}^N از توزیع یکنواخت تبعیت نمایند. بر اساس نتایجی که آن‌ها در مقاله خود گزارش نمودند، این مدل از سرعت کمتری نسبت به مدلی که در فضای اندازه-شکل تعریف می‌شود، برخوردار است. در این نوع مدل‌بندی نیز مشکلات مربوط به همگرایی و وابستگی به نقطه شروع وجود دارد. از این‌رو در بخش ۳ روش‌های بهبود این الگوریتم‌ها معرفی خواهد شد و در بخش ۴ با استفاده از داده‌های واقعی عملکردشان مورد مقایسه قرار خواهد گرفت.

۳ بهبود مدل های انطباق

فرض کنید ماتریس های Z و ζ نشان های مربوط به پیکربندی X و μ باشند، به طوری که ζ به صورت ثابت در نظر گرفته شود و توزیع های X و Z به صورت شرطی از یکدیگر مستقل باشند. اگر θ پارامتر (های) توزیع مشخصه Z باشد، می توان تابع درستنمایی مدل جدید را به صورت

$$L(X, Z|\Lambda, \mu, \tau, \zeta, \theta) = L(X|\Lambda, \mu, \tau)L(Z|\zeta, \theta)$$

نوشت. اولین بار درآیدن و همکاران (۲۰۰۷) این مدل بندی را بر اساس توزیع نرمال برای اضافه کردن اطلاعات بار جزئی و شعاع و اندروال به مدل احتمال در مسئله انطباق استروئیدها استفاده نمودند. در عمل استفاده از این اطلاعات سبب بهبود نسبی سرعت همگرایی الگوریتم *MCMC* در انطباق استروئیدها می شود. در مسئله انطباق پروتئین نیز می توان با استفاده از همین روش خواص شیمیایی و فیزیکی پروتئین ها مانند اندازه، قطبیت، بار جزئی، هیدروفوبیسیته را به مدل سازی احتمالاتی اضافه نمود. اما از آنجا که ساختار پروتئین در مقایسه با ساختار استروئیدها بسیار بزرگتر و پیچیده تر است، در نظر گرفتن یک توزیع احتمال کلی برای نشان های نقاط، بهبودی را در روند الگوریتم ایجاد نمی کند.

اکثر الگوریتم های قطعی از نتایج انطباق ساختار اول، دوم و خواص شیمیایی و فیزیکی پروتئین ها در انطباق ساختار سوم بهره می برند. در روش های احتمالاتی نیز می توان از این اطلاعات به منظور افزایش سرعت و دقت الگوریتم استفاده کرد. بر اساس همین انگیزه، فالیزی و همکاران (۲۰۱۴) و نجیبی و همکاران (۲۰۱۵) با تعمیم مدل های گرین و ماردیا (۲۰۰۶) و درآیدن و همکاران (۲۰۰۷)، مدل بندهای جدیدی را بر اساس انطباق اطلاعات ساختار اول سوم به صورت همزمان معرفی کرده اند. آن ها توزیع پیشین مناسبی را بر مبنای انطباق ساختار اول با در نظر گرفتن گپ های بین توالی اسید آمینه ها ارائه داده اند که در عمل سرعت الگوریتم را در رسیدن به همگرایی مطلوب افزایش می دهد. ایده اصلی این روش ها انطباق موضعی نقاط به جای انطباق سراسری آن ها بر اساس اطلاعات دیگری به جز اطلاعات ساختار سوم است. به عنوان مثال در روش نجیبی و همکاران (۲۰۱۵)،

ابتدا یک جورسازی موضعی بر اساس هرم‌های دلاونی ایجاد شده برای هر مشخصه پروتئین انجام می‌شود و سپس این اطلاعات در قالب توزیع پیشین یا توزیع نامزد در الگوریتم متروپولیس هستینگز، استفاده می‌شود. نتایج شبیه‌سازی استفاده از این روش‌ها در انطباق پروتئین‌ها نشان داده است که سرعت و دقت الگوریتم بهبود می‌یابد. دلیل اصلی این بهبود را می‌توان وجود اطلاعات ساختاری موضعی، میزان مشابهت ساختار اول بر اساس هر یک از ماتریس‌های جانشین و اطلاعات توالی اسیدهای آمینه در این مدل‌ها دانست.

۴ انطباق احتمالاتی پروتئین‌ها

هدف از این بخش انطباق زوجی شکل سه‌بعدی یک مجموعه از پروتئین‌ها شامل ۵ زوج پروتئین است. این مجموعه به تصادف از بین ۱۶ زوج پروتئینی که اولین بار توسط اریتز و همکاران (۲۰۰۲) مورد استفاده قرار گرفته است، انتخاب شده است. داده‌های مربوط به این پروتئین‌ها از بانک اطلاعاتی پروتئین به آدرس <http://www.pdb.org> قابل دسترسی هستند. برای آشنایی با نحوه توالی اسیدهای آمینه و ساختار دوم آنها، نمودار سه‌بعدی دو پروتئین 1ACX با ۱۰۸ اسید آمینه و 1COB با ۱۵۱ اسید آمینه در شکل ۱ رسم شده است. به منظور مقایسه نتایج این مجموعه دو حالت متفاوت برای مدل احتمالاتی درآیدن و همکاران (۲۰۰۷) در نظر گرفته شده است. (الف) از اطلاعات توالی و نوع اسید آمینه در انطباق استفاده نشود (ب) بر اساس نجیبی و همکاران (۲۰۱۵)، اطلاعات نوع اسید آمینه و توالی آنها به کمک مثلث‌سازی دلاونی و میانگین وزنی تابع مشابهت ساختار سوم و نوع اسید آمینه به کمک ماتریس بلوکی جانشین^{۱۰} (BLOSUM62) به مدل اضافه شود. در تمامی حالت‌ها به منظور شروع الگوریتم ۶ نقطه به تصادف انتخاب و بعد از ۱۰۰۰۰۰ تکرار، همگرایی الگوریتم بررسی و در صورت عدم همگرایی، الگوریتم مجدداً اجراء شده است. این روش تا همگرایی لازم تکرار و پس از تایید همگرایی، الگوریتم تا ۴۰۰۰۰۰ تکرار ادامه داده و پارامترهای مدل بر اساس نمونه‌های مستقل

^{۱۰} BLOcks SUBstitution Matrix

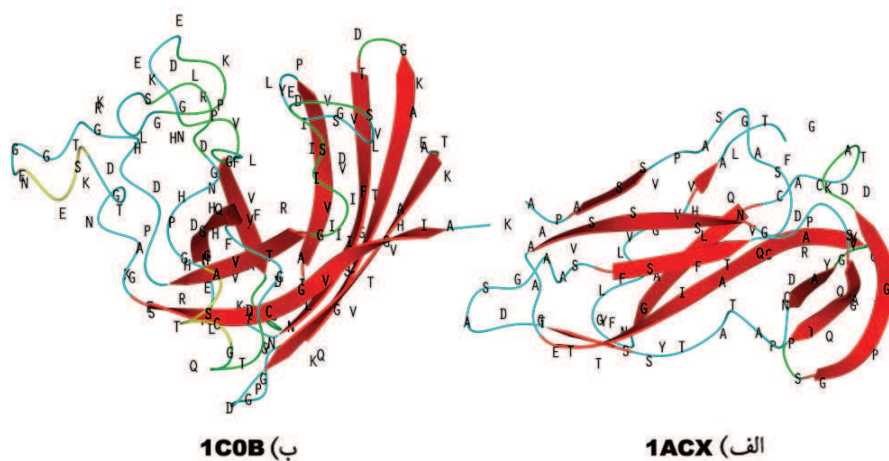
جدول ۱: نتایج مدل انطباق احتمالاتی در فضای اندازه-شکل و روش قطعی CE در ۵ جفت پروتئین به همراه بازه اعتبار پسین ۹۵ درصد پارامترها

الگوریتم CE		(ب)		(الف)				
<i>RMSD</i>	<i>p</i>	<i>RMSD</i>	\hat{p}	<i>RMSD</i>	\hat{p}	ℓ	k	پروتئین
۴/۵۹	۵۴	۳/۶۱	۵۶	۳/۷۸	۵۷	۱۸۸	۸۷	1ABA-1DSB
		(۳/۴۵, ۴/۲۱)	(۵۲, ۵۹)	(۳/۴۵, ۴/۲۱)	(۵۲, ۵۹)			
۴/۰۵	۹۳	۴/۰۹	۹۱	۳/۹۸	۸۹	۱۵۱	۱۰۸	1ACX-1COB
		(۳/۸۴, ۴/۲۵)	(۸۴, ۹۵)	(۳/۸۷, ۴/۳۱)	(۸۳, ۹۴)			
۳/۱۳	۴۹	۳/۲۱	۵۴	۳/۱۶	۵۱	۱۹۴	۵۶	1PGB-5TSS
		(۲/۵۴, ۳/۸۰)	(۴۷, ۵۸)	(۲/۴۱, ۳/۶۷)	(۴۴, ۵۵)			
۴/۱۰	۱۱۵	۴/۰۳	۱۱۱	۴/۱۳	۱۱۰	۱۸۵	۱۵۲	1TNF-1BMV
		(۳/۷۸, ۴/۳۱)	(۱۰۲, ۱۱۶)	(۳/۷۳, ۴/۲۴)	(۱۰۱, ۱۱۵)			
۳/۵۹	۸۵	۲/۹۱	۷۹	۲/۹۸	۸۱	۱۰۶	۱۵۴	2TMV-256B
		(۲/۵۶, ۳/۴۴)	(۷۵, ۸۵)	(۲/۶۵, ۳/۴۴)	(۷۶, ۸۵)			

برآورد می‌شود. به منظور مقایسه نتایج الگوریتم‌های احتمالاتی با یکدیگر، برآورد تعداد نقاط جور شده (\hat{p}) و جذر میانگین مربع انحرافات^{۱۱} ($RMSD$) به همراه بازه اعتبار پسین ۹۵ درصد آنها و تعداد نقاط جور شده (p) و $RMSD$ در الگوریتم CE (شیندیالو و بورنه، ۱۹۹۸) در جدول ۱ گزارش شده است.

در اجرای مدل (الف) برای هر یک از ۵ زوج پروتئین، الگوریتم به‌طور میانگین پس از دهمین تکرار (حداقل ۶ و حداکثر ۱۸ تکرار) به همگرایی لازم دست یافت. در حالی که در اجرای مدل (ب) که در آن از اطلاعات اسیدهای آمینه و توالی آنها در مدل‌سازی استفاده شده است، الگوریتم به‌طور میانگین در سومین تکرار (حداقل ۱ و حداکثر ۵ تکرار) به همگرایی لازم دست یافت. تعداد نقاط جور شده در روش احتمالاتی و روش CE در هیچ کدام از پروتئین‌ها به‌طور ۱۰۰ درصد یکسان نبوده است. همان‌طور که ملاحظه می‌شود تعداد نقاط جور شده در روش CE در دو زوج پروتئین‌های 1ABA-1DSB و 2TMV-256B در بازه اعتبار پسین مدل‌های احتمالاتی قرار گرفته است، اما مقدار $RMSD$ هیچ کدام از آنها در بازه اعتبار پسین متناظر

^{۱۱} Root Mean Square Deviation



شکل ۱: دو پروتئین و توالی کربن‌های آلفا به همراه اطلاعات نوع اسید آمینه و نوع ساختار دوم در فضای سه بعدی الف) پروتئین 1ACX با ۱۰۸ اسید آمینه (ب) پروتئین 1C0B با ۱۵۱ اسید آمینه، صفحات مشکلی رنگ نشان دهنده‌ی ساختار دوم صفحات موازی بتا و حروف نمایانگر نوع اسید آمینه است.

قرار نگرفته است. بنابراین می‌توان گفت که در این دو زوج پروتئین با احتمال ۹۵ درصد الگوریتم *CE* به خوبی عمل نکرده است. از آن‌جا که در داده‌های واقعی جواب درست و حقیقی وجود ندارد و مقایسه الگوریتم‌ها به سادگی انجام نمی‌شود، نمی‌توان هیچ کدام از این روش‌ها را به طور قطعی پذیرش یا رد کرد. اما با توجه به نتایج فوق، می‌توان نتیجه گرفت که با نادیده گرفتن زمان اجرای روش‌های احتمالاتی، نتایج انطباق آنها منطقی و قابل مقایسه و رقابت با روش‌های قطعی است. از طرفی دیگر، بازه اعتبار پسین و دیگر ابزارهای آماری می‌توانند در مقایسه نتایج الگوریتم‌ها استفاده شوند. در انتها برای ارائه نمایشی از نتایج نهایی انطباق، تصویر پروتئین‌های منطبق شده 1ACX-1C0B بر یکدیگر بر اساس ۹۱ نقطه جور شده که از برآورد پارامترها توسط مدل احتمال در فضای اندازه-شکل به همراه اطلاعات اسیدهای آمینه و توالی آنها به دست آمده است، در شکل ۲ رسم شده است.



شکل ۲: توالی کربن‌های آلفا پروتئین IACX (رنگ مشکی) زمانی که بر پروتئین 1C0B (رنگ خاکستری) با ۹۱ نقطه جور شده است ($RMSD = ۴/۰۹$)

بحث و نتیجه گیری

در مقاله حاضر روش‌های احتمالاتی مسئله انطباق ساختاری ارائه و با تشکیل تابع درست‌نمایی، تعیین پیشین‌های مناسب استنباط پسینی انجام شد. با این حال، لزوم تعمیم روش‌های آماری در این حوزه و ارائه راهکارهای مناسب برای افزایش سرعت و دقت الگوریتم‌های پیشنهادی ضروری است. یکی از راه‌حل‌های مقابله با تکرار بسیار زیاد الگوریتم‌ها، استفاده از اطلاعات شیمیایی و فیزیکی پروتئین‌ها مانند خواص هیدروفوبیسیته، اندازه، قطبیت و ... در انطباق موضعی پروتئین‌ها و در نهایت انطباق سراسری آنها است. علاوه بر این، استفاده از نمایش زاویه‌ای بجای نمایش مختصاتی پروتئین و بهره‌گیری از آمار جهت‌دار ممکن است موجب بهبود بهتر روش‌ها شود.

تقدیر و تشکر

نویسندگان مقاله بر خود لازم دانسته از نظرات ارزشمند سردبیر، هیئت تحریریه و سه داور محترم مجله علوم آماری که سبب ارتقاء کیفیت مطالب ارائه شده در این مقاله شدند، تشکر و قدردانی نمایند.

مراجع

- Dryden, I., Hirst, J. and Melville, J. (2007), Statistical Analysis of Unlabeled Point Sets: Comparing Molecules in Chemoinformatics, *Biometrics*, **63**, 237-251.
- Dryden, I. L. and Mardia K. V. (1998), *Statistical Shape Analysis*, John Wiley, Chichester.
- Fallaize, C., Green, P., Mardia, K. and Barber, S. (2014), Bayesian Protein Sequence and Structure Alignment, *arXiv preprint*, *arXiv:1404.1556*.
- Green P. J. and Mardia K. V. (2006), Bayesian Alignment Using Hierarchical Models, with Applications in Protein Bioinformatics, *Biometrika*, **93**, 235-254.
- Kendall, D. G., Barden, D., Carne, T. K. and Le, H. (1999), *Shape and Shape Theory*, John Wiley, Chichester.
- Kent, J. T. and Mardia, K. V. (2001), Shape, Tangent Projections and Bilateral Symmetry, *Biometrika*, **88**, 469-485.
- Kent, J. T., Mardia, K. V. and Taylor C. C. (2004), Matching Problems for Unlabelled Configurations, *Bioinformatics, Images, and Wavelets, Proceedings of LASR 2004*, Ed. R. G. Aykroyd, S. Barber and K. V. Mardia, University of Leeds, Leeds, 33-36.
- Kenobi, K. and Dryden, I. L. (2012), Bayesian Matching of Unlabelled Point Sets Using Procrustes and Configuration Models, *Bayesian Analysis*, **7**, 1-20.

- Le, H. (1988), *Shape Theory in Flat and Curved Spaces, and Shape Densities with Uniform Generators*, Ph.D thesis, University of Cambridge.
- Mardia, K. V., Fallaize, C. J., Barber, S., Jackson, R. M. and Theobald, D. L. (2013), Bayesian Alignment of Similarity Shapes, *The Annals of Applied Statistics*, **7**, 989-1009.
- Najibi, S. M., Faghihi, M. R., Golalizadeh, M. and Arab, S. S. (2015), Bayesian Alignment of Proteins via Delaunay Tetrahedralization, *Journal of Applied Statistics*, **42**, 1064-1079.
- Ortiz, A. R., Strauss, C. E. M. and Olmea, O. (2002), Mammoth (Matching Molecular Models Obtained from Theory): An Automated Method for Model Comparison, *Protein Science*, **11**, 2606-2621.
- Shindyalov, I. N. and Bourne, P. E. (1998), Protein Structure Alignment by Incremental Combinatorial Extension (CE) of the Optimal Path, *Protein Engineering Journal*, **11**, 739-747.
- Shindyalov, I. N. and Bourne, P. E. (2000), An Alternative View of Protein Fold Space, *Proteins: Structure, Function, and Bioinformatics*, **38**, 247-260.
- Taylor, C. C., Mardia, K. V. and Kent, J. T. (2003), Matching Unlabelled Configurations Using the EM Algorithm, *Stochastic Geometry, Biological Structure and Images, Proceedings of LASR 2003*, Ed. Aykroyd, R. G., Mardia, K. V. and Langdon, M. J., University of Leeds, 19-21.