



روشی

جدید برای نسبت دادن ساختار دوم پروتئین

بر اساس روابط ریاضیاتی بین مختصات سه بعدی C_{α} هاحسینی سیدرزگار^{1و2*}، صادقی مهدی^{3**}، حبیبی مهناز^{2و4}، اصلاح چي چنگیز⁴

1- گروه بیوتکنولوژی، پردیس علوم، دانشگاه تهران 2- هسته بیوانفورماتیک، پژوهشگاه علوم کامپیوتر، پژوهشگاه دانش های بنیادی 3- پژوهشگاه ملی مهندسی ژنتیک و زیست فناوری، تهران 4- دانشکده علوم ریاضی، دانشگاه شهید بهشتی

چکیده فارسی: نسبت دادن ساختار دوم پروتئین، گام اصلی برای آنالیز و مدلسازی ساختمان پروتئین است. به همین دلیل روشهایی با معیارهای مختلف به این منظور بوجود آمده اند. ما روش جدیدی را برای نسبت دادن ساختار دوم به پروتئین با پیدا کردن چهار رابطه ریاضیاتی بین مختصات سه بعدی C_{α} های اسید آمینه های متوالی، معرفی کرده ایم که انطباق با این روابط را ملاک نسبت دادن ساختار دوم به اسید آمینه ها قرار می دهیم.

کلمات کلیدی: ساختار دوم پروتئین، مختصات سه بعدی C_{α} ، رابطه های ریاضیاتی

مقدمه: ساختارهای دوم، ساختارهایی در پروتئین هستند که به دلیل

برقراری پیوند هیدروژنی بین اسید آمینه ها بر اساس یک الگوی خاص، به

وجود می آیند. این ساختارها را در سه دسته کلی طبقه بندی کرده

اند: **helix.1** ها (که شامل مارپیچ های α ، 3_{10} و π می باشد.) **2. β -strand** ها **3**

coil. ها که معرف ساختارهای نامنظم موجود در پروتئین است. ساختارهای دوم

به این دلیل که توصیف ساده و شهودی از ساختمان سه بعدی پروتئین

فراهم می کنند، در زیست شناسی ساختمانی کاربردهای فراوانی پیدا کرده

اند. از جمله در طبقه بندی پروتئین ها، ردیف بندی بهتر توالی ها، ردیف

بندی ساختمان پروتئینها و مدلسازی ساختمان سه بعدی به کار می

روند. بنابراین وجود روشهایی که بتواند ساختمان دوم را از مختصات سه

بعدی آنها تعیین کند، ضروری است. اوایل کریستالوگرافها با چشم این کار

را انجام می دادند اما با افزایش تعداد پروتئینهای تعیین ساختار

شده، لزوم وجود روش هایی که بصورت اتوماتیک این کار را انجام دهد در

اواسط دهه 70 احساس شد و روشهای مختلفی برای این کار بوجود آمدند. از

جمله این روشها **DSSP**¹ است که الگوی تکرار پیوندهای هیدروژنی را با

حساب کردن انرژی پیوند زنجیره اصلی پیدا می کند و بر اساس آن ساختار

دوم نسبت می دهد. **STRIDE**² که مشابه **DSSP** است با این تفاوت که نحوه

محاسبه پیوند هیدروژنی اش متفاوت بوده و معیار زوایای پیوند را نیز به

میان می آورد. علاوه بر این دو، روش های هندسی فراوانی بوجود آمده

اند، مثل **DEFINE** و **P-CURVE**، **SECSTER**، **VOTAP** و... ما روشی را

بدین منظور تحت عنوان **MATHREL** (**MATH**mathical **REL**ation based assignment) ارائه

کرده ایم که اساس این روش بر 4 رابطه ریاضیاتی که بین مختصات C_{α} اسید

آمینه های مجاور پیدا کرده ایم، استوار است که هر کدام از این رابطه

ها به یکی از ساختارهای منظم (**β -strand, π -helix, 3_{10} -helix, α -helix**) تعلق دارد و در

نهایت با استفاده از این رابطه ها ساختارهای دوم را به اسید آمینه های

موجود در پروتئین نسبت می دهیم. بمنظور بررسی اعتبار روش خود، آن را

با روش های رایج **STRIDE**، **DSSP** و فایل های **PDB**³ مقایسه کرده ایم و نشان

دادیم که با آن ها انطباق خوبی دارد.

روش: در این قسمت مراحل بدست آوردن 4 رابطه ریاضیاتی مذکور را معرفی

می کنیم.

با استفاده از این خاصیت که یک مارپیچ در امتداد هر یک از محورهای

مختصات دارای 2 ویژگی نوسان و امتداد خطی است. معادله پارامتری مارپیچ

را در حالت کلی معرفی کرده ایم:

$$g(t)=f(x(t), y(t), z(t)),$$

$$x(t)=a_x \sin(\omega t+ \theta_0)+ b_x t+ c_x$$

$$a_x= R \sin \varphi_x, b_x= R \cos \varphi_x,$$



$$y(t) = a_y \sin(\omega t + \pi/2 + \theta_0) + b_y t + c_y \quad a_y = R \sin \phi_y, b_y = R \cos \phi_y,$$

$$z(t) = a_z \sin(\omega t + \pi + \theta_0) + b_z t + c_z \quad a_z = R \sin \phi_z, b_z = R \cos \phi_z$$

که R شعاع مقطع مارپیچ، ϕ زاویه بین محور مارپیچ و محور مختصات متناظر با آن، ω چرخش زاویه ای نوسان، B ارتفاع هر دور مارپیچ، θ_0 زاویه اولیه نوسان بوده و c_x, c_y, c_z ثابت های اولیه هستند. به منظور بدست آوردن روابط کلی که برای همه مارپیچ ها صادق باشد، پارامترهایی را که به یک مارپیچ خاص تعلق دارند، مثل a_x, b_x, c_x را طی مراحل زیر حذف می کنیم.

فرض کنیم $c(i)$ و $0 \leq i \leq 4$ ، پنج نقطه متوالی در یک مارپیچ باشند آن گاه بر اساس رابطه (1) داریم:

$$c(i) : x_i = a_x \sin(\omega i + \theta_0) + b_x i + c_x \quad (2)$$

به منظور حذف b_x و c_x ، $A(j)$ و $j=0$ را تعریف می کنیم، به طوری که

$$A(j) = x_{(j+1)} - 2x_{(j)} + x_{(j-1)} \quad (3)$$

حال برای حذف a_x ، $A(j)$ را بر $A(j+1)$ تقسیم می کنیم، بنابراین داریم:

$$\frac{A_x(j)}{A_x(j+1)} = \frac{\lambda(j) \cos \theta_0 + \gamma(j) \sin \theta_0}{\lambda(j+1) \cos \theta_0 + \gamma(j+1) \sin \theta_0} \quad (4)$$

که در آن

$$\lambda(j) = \sin((j+2)\omega) - 2\sin((j+1)\omega) + \sin(j\omega)$$

$$\gamma(j) = \cos((j+2)\omega) - 2\cos((j+1)\omega) + \cos(j\omega) \quad (5)$$

در مرحله آخر منظور حذف θ_0 ، رابطه نهایی را تعریف می کنیم:

$$\frac{\frac{A_x(j)}{A_x(j+1)} - \frac{\lambda(j)}{\lambda(j+1)}}{\frac{A_x(j+1)}{A_x(j+2)} - \frac{\lambda(j+1)}{\lambda(j+2)}} = \frac{\lambda(j+1)\gamma(j) - \lambda(j)\gamma(j+1)}{\lambda(j+1)\gamma(j+2) - \lambda(j+2)\gamma(j+1)} \quad (6)$$

که این رابطه مستقل از پارامترهای مذکور است و فقط به λ و γ بستگی

دارد که آن ها طبق (5) فقط به ω بستگی دارند و $\omega = \frac{2\pi}{T}$ که T دوره تناوب است.

این دوره تناوب برای مارپیچ های موجود در پروتئین بصورت تعداد C_α ها در یک دور تعریف می شود که آن را در جدول [1] نمایش داده ایم. که با استفاده از ω می توانیم λ و γ را برای هر یک از ساختارها حساب کرده و با استفاده از معادله های (3) و (6) معادله مخصوص هر یک از این مارپیچ ها را بدست می آوریم. برای 5 نقطه متوالی در یک توالی پروتئینی با N اسید آمینه $1 \leq k \leq N-4$ ، 3 رابطه زیر را برای هر یک از مارپیچ ها داریم. که برای مارپیچ های α, π و m به ترتیب 2.88، 3.51- و 1 و n به ترتیب 2.53-، 3.23 و 0 است.

$$f(x) = \frac{\left(\frac{x_{k+2} - 2x_{k+1} + x_k}{x_{k+3} - 2x_{k+2} + x_{k+1}}\right) + m}{\left(\frac{x_{k+4} - 2x_{k+3} + x_{k+2}}{x_{k+3} - 2x_{k+2} + x_{k+1}}\right) + n} = -1 \quad (7)$$



در واقع می شود β -strand را به عنوان مارپیچی در نظر گرفت که در هر دور، 2 تا C_{α} دارد و بنابراین $T=2$ و $\omega=\pi$ ، پس با ساده تر کردن رابطه، برای هر 4 نقطه متوالی در یک توالی با N اسید آمینه $1 \leq k \leq N-3$ داریم:

$$f_{\beta}(x_k) = \frac{x_{(k+3)} - x_{(k+1)}}{x_{(k+2)} - x_{(k)}} = 1. \quad (8)$$

الگوریتم: در ابتدا یک پنجره متحرک 5 تایی که در هر حرکت یک آمینو اسید جلو می رود را در نظر می گیریم، سپس مقادیر $(f(x), f(y), f(z))$ را براساس معادله (7) برای هر 3 نوع مارپیچ و هر کدام از پنجره ها محاسبه می کنیم. و بعد انحراف هر پنجره را از هر کدام از مارپیچ ها بدست می آوریم. و این مراحل را برای β -strand با تعریف کردن پنجره متحرک 4 تایی و استفاده از رابطه (8) تکرار می کنیم و در نهایت مقدار انحراف هر اسید آمینه را از هر کدام از ساختارها با استفاده از انحراف پنجره هایی که در آن قرار می گیرد حساب می کنیم سپس بر اساس آن، ساختار دوم را به تک اسید آمینه ها نسبت می دهیم.

جدول [1]: خاصیت های

	α -helix	π -helix	3_{10} -helix
T	3.6	4.4	3
ω	$2\pi/3.6$	$2\pi/4.4$	$2\pi/3$

تناوبی مارپیچ ها

نتایج و بحث: روش ما، MATHREL چون فقط از مختصات C_{α} ، برای نسبت دادن ساختار دوم استفاده می کند، این باعث افزایش سرعت، سادگی استفاده و بازدهی بالای آن شده است. علاوه بر اینها برای اعتباربخشی روشمان آن را با PDB, STRIDE, DSSP مقایسه کرده ایم. برای این مقایسه یک مجموعه شامل 1918 پروتئین با شباهت کمتر از 90% و قدرت تفکیک بهتر از 2A را مورد استفاده قرار دادیم و نتایج موجود در جدول [2] بدست آمد.

بر اساس جدول [2]، درصد توافق مارپیچها بین روش ما و هر 3 روش بالاتر از 90% است که این بیانگر آن است که در پیدا کردن مارپیچها، MATHREL با این روشهای معتبر، به خوبی موافقت دارد و با استناد به آنها مارپیچ ها را خوب تشخیص می دهد. از طرف دیگر درصد توافق در پیدا کردن β -strand کمی پایین تر است. دلیل اصلی این مساله این است که روشهای براساس پیوند هیدروژنی، فقط قادر به تشخیص β -strand هایی هستند که در پیوند هیدروژنی شرکت کرده اند، است و نمی توانند β -bridge ها را پیدا کنند. در حالیکه روش ما آن ها را تشخیص می دهد. وی در کل با وجود تفاوت زیاد در الگوریتم MATHREL با این روشها، توافق رضایت بخشی با این روش های معتبر بدست آمده است که اعتبار روش ما را تایید می کنند. برای مقایسه MATHREL با روش های DSSP, STRIDE و فایل های موجود در PDB ما از معیارهای شناخته شده از قبیل درصد توافق ($\% Sn$ (Sensitivity))

$$Sp \text{ (Specificity)} \text{ استفاده کرده ایم، } Sp = \frac{TP}{TP + FP} \text{ و } Sn = \frac{TP}{TP + FN}$$

که در آن TP (True Positive), TN (True Negative), FN (False Negative) و FP (False Positive)

جدول [2]: نتایج مقایسه MATHREL با روش های دیگر:

ساختار Sp % Sn کل FN FP TN TP روش های مقایسه شده



MATHREL_PDB	165145	298456	11468	39541	514610	90.09	81.20	92.81	Helix
MATHREL-DSSP	152995	316351	23618	21646	514610	91.20	87.61	86.63	Helix
MATHREL- STRIDE	158452	316084	18161	21913	514610	92.21	87.85	87.85	Helix
MATHREL-PDB	63369	350001	46102	55138	514610	80.32	53.47	57.89	β -Strand
MATHREL-DSSP	68334	349306	41137	55833	514610	81.16	55.03	62.42	β -Strand
MATHREL-STRIDE	70997	348882	38474	56257	514610	81.59	55.79	64.85	β -Strand

تشکر و قدردانی: بخشی از هزینه این کار از محل طرح شماره CS 1385-1-02 پژوهشگاه دانش های بنیادی (IPM) تامین شده است.

منابع:

- 1 Kabsch,W., Sander,C. (2001) Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical Features, Biopolymers, 22(12), 2577-637
- 2 Frishman,D., Argos,P. (1995) Knowledge-based protein secondary structure assignment, proteins, 23(4), 566-579.
- 3 Sussman JL, Lin DW, Jiang JS, Manning NO, Pirlusky J et al(1998). Protein Data Bank (PDB): Database of three-dimensional structural information of biological macromolecules, Acta Crystallographica Section D-Biological Crystallography 54,1078-84

Protein Secondary Structure Assignment Based on Mathematical Relations Between $C\alpha$ Three Dimensional Coordinates

Sayed Rzgar Hosseini^{1,2,*}, Mehdi Sadeghi^{3**}, Mahnaz Habibi^{2,4}, Changiz Eslahchi⁴

1- Biotechnology group, College of Science, University of Tehran

2- School of Computer Science, Institute for Studies in Theoretical Physics and Mathematics

3- National Institute of Genetic Engineering and Biotechnology

4- Faculty of Mathematics, Shahid Beheshti University

Abstract:

The automatic assignment of protein secondary structure from three-dimensional atomic coordinate of proteins is an essential step for the analysis and modeling of protein structures. So different methods regarding different criteria have been designed to perform this task. We introduce a new method for protein secondary structure assignment based solely on $C\alpha$ coordinates. We have found four mathematical relations between three-dimensional coordinates of consecutive residues, each of which applies to one of the four regular secondary structure categories (α -helix, 310-helix, π -helix and \downarrow -strand). Our algorithm calculates deviation of the $C\alpha$ coordinate of each residue from each of these relations. Then based on a defined cutoff for deviation from each relation, it assigns secondary structures to all residues of a protein.

Keywords: protein secondary structure assignment, $C\alpha$ three-dimensional coordinate, mathematical relations



جمعین ہمیش ملی بیوتکنولوژی جمہوری اسلامی ایران
3-5 آذر ماہ 1386، سالن اجلاس سران
The 5th National Biotechnology Congress of Iran
24-26 Nov, 2007, Summit Meeting Conference Hall, Tehran- Iran

