



## ترکیبی از مدل های مارکوفی پنهان و شبکه های عصبی برای تعیین ساختار دوم پروتئین ها

- پزشک حمید<sup>1</sup> و<sup>2</sup>\*\*\* صادقی مهدی<sup>3</sup> اصلاحی چنگیز<sup>4</sup> ملک پورسید امیر<sup>1</sup> و<sup>2</sup>  
1 قطب بیومتمتیک و بخش آمار در دانشکده ریاضی، آمار و علوم کامپیوتر،  
پردیس علوم دانشگاه تهران  
2 هسته تحقیقاتی بیو انفورماتیک، پژوهشکده کامپیوتر، پژوهشگاه داده  
نشانی بنیادی (IPM)  
3 پژوهشگاه ملی مهندسی ژنتیک و زیست فناوری  
4 دانشکده علوم ریاضی دانشگاه شهید بهشتی

### چکیده :

پروتئین ها از توالی اسید های آمینه ساخته شده اند و برای هر یک از اسید های آمینه یکی از سه ساختار Helix (H)، Strand (S) یا Coil (C) در نظر گرفته می شود که ساختار دوم پروتئین ها را می سازد. این ساختار برای بسیاری از توالی های پروتئینی نا مشخص است. پیش بینی ساختار دوم پروتئین ها یکی از مراحل مهم در تعیین ساختمان سه بعدی و عملکرد پروتئین ها به شمار می رود. شبکه های عصبی معمولاً به طور مستقیم برای محاسبه احتمال های نگاشت در مدل های مارکوفی پنهان به کار گرفته می شوند. در این مقاله از ساختار ترکیبی جدیدی از مدل های مارکوفی پنهان و شبکه های عصبی استفاده شده است به طوری که ساختار پیش بینی شده توسط شبکه های عصبی را به عنوان یک پیش بینی اولیه برای محاسبه احتمال نگاشت اسید های آمینه در مدل مارکوفی پنهان به کار گرفته ایم و پیش بینی جدیدی بر اساس مدل های مارکوفی پنهان انجام داده ایم. با استفاده از این ساختار ترکیبی دقت پیش بینی مدل های مارکوفی پنهان را در تعیین ساختار دوم پروتئین ها افزایش داده ایم.

### واژه های کلیدی:

مدل های مارکوفی پنهان، شبکه های عصبی، ساختار دوم پروتئین، روش های بی زی

### 1- مقدمه :

امروزه روش های متعددی برای پیش بینی ساختار دوم پروتئین بر اساس توالی اسید های آمینه پیشنهاد شده است. پیش بینی ساختار دوم پروتئین از منظر مدل های مولد در (1) و (2) بررسی شده است که یک مدل مولد آماری را بر اساس مدل های مارکوفی پنهان گسترش یافته<sup>1</sup> توسعه داده است و از یک دوره به طول متغیر برای وضعیت های آن استفاده می شود چنین مدلی اجازه می دهد در هر وضعیت پنهان طول متغیری از مشاهدات نگاشته شود و این به مدل نیمه مارکوفی قطاعی<sup>2</sup> معروف شده است یکی از فواید استفاده از چنین روش های احتمالی این است که می توان منابع متفاوتی از اطلاعات موجود در توالی را با استفاده از توزیع احتمال توام ساختار-توالی و بر اساس ساختار های قطاعی وارد مدل کرد. در (3)، (4) با استفاده از اطلاعات تکاملی پیش بینی بهتری برای SSMM بدست آمده است.

در این مقاله مدل جدیدی برای بهبود دقت SSMM ها به کار گرفته شده است که یک ساختار ترکیبی از شبکه های عصبی و SSMM ها را بدون در نظر گرفتن اطلاعات تکاملی به کار می گیرد. ابتدا اسید آمینه ها را به صورت پنجره هایی به شبکه عصبی داده ایم و سپس پیش بینی های انجام شده توسط شبکه عصبی را برای تخمین

<sup>1</sup> Generalized Hidden Markov Models (GHMM)  
<sup>2</sup> Segmental Semi-Markov Models (SSMM)



پارامتر های نگاشت در مدل مارکوفی پنهان و بالا بردن دقت پیش بینی بدست آمده در مرحله اول توسط شبکه عصبی، به کار گرفته ایم.

## 2- شبکه های عصبی :

در این مقاله شبکه های عصبی برای انجام یک پیش بینی اولیه مورد استفاده قرار گرفته است به این صورت که توالی اسید های آمینه را به صورت پنجره های 9 تایی به عنوان ورودی شبکه عصبی در نظر گرفته ایم و یک پیش بینی اولیه را برای اسید آمینه ای که در مرکز این پنجره قرار دارد بدست آورده ایم. برای کاهش تعداد نرون های ورودی شبکه عصبی از یک بردار 6 تایی برای کد کردن اسید های آمینه استفاده کرده ایم. همچنین خروجی شبکه عصبی شامل سه نرون است. شبکه عصبی مورد استفاده یک شبکه پیشرو دو لایه است که خطا های بدست آمده در هر مرحله را در شبکه گسترش می دهد<sup>3</sup> همچنین 20 نرون برای لایه پنهان در نظر گرفته شده است.

## 3- مدل های مارکوفی پنهان :

ساختار دوم پروتئین می تواند توسط یک سه تایی (m,e,T) تعریف شود که m نشان دهنده تعداد قطاع<sup>4</sup> های توالی است. e نقاط انتهایی قطاع ها و T ساختارقطاع ها را نشان می دهد. در شکل (1) یک مدل گرافیکی از ساختار دوم پروتئین برای T=(C,S,C,H,...) و e=(3,5,7,11,...) نشان داده شده است. هر مشاهده O<sub>i</sub> شامل اسید آمینه R<sub>i</sub> و ساختار پیش بینی شده توسط شبکه عصبی، P<sub>i</sub> است همچنین کل مشاهدات را با O نشان می دهیم. برای محاسبه توزیع پسین (m,e,T) به شرط توالی مشاهدات، محاسبه توزیع پیشین برای متغیر های (m,e,T) که ساختار دوم پروتئین را مشخص می کند ضروری است و آن را به صورت زیر نویسیم :

$$P(m, e, T) = P(m)P(e, T | m) = P(m) \prod_{i=1}^m P(e_i | e_{i-1}, T_i) P(T_i | T_{i-1})$$

P(m) احتمال پیشین مشاهده m قطاع را در یک توالی نشان می دهد. P(e<sub>i</sub> | e<sub>i-1</sub>, T<sub>i</sub>) به ما اجازه می دهد که طول قطاع ها را به الگو در آوریم و آن را به شکل زیر بیان می کنیم :

$$P(e_i | e_{i-1}, T_i) = P(e_i - e_{i-1} | T_i)$$

که احتمال نگاشت یک مشاهده به طول e<sub>i</sub>-e<sub>i-1</sub> رادر وضعیت T<sub>i</sub> نشان می دهد. چون در مرحله استنباط از m استفاده ای نمی کنیم یک توزیع یکنواخت برای آن در نظر می گیریم. برای محاسبه P(m,e,T|O) باید مقدار

P(O | m, e, T) را نیز محاسبه کنیم که به صورت زیر تعریف می شود :

$$P(O | m, e, T) = \prod_{i=1}^m P(S_i | S_{i-1}, T_i)$$



$S_i$  مشاهدات قطاع  $i$  ام را نشان می دهد  $S_{-i}$  مشاهدات قطاع های قبلی را نشان می دهد. برای  $P(S_i | S_{-i}, T_i)$  داریم :

$$P(S_i | S_{-i}, T_i) = \prod_{k=e_{i-1}+1}^{e_i} P(O_k | O_{[1:k-1]}, T_i) = \prod_{k=e_{i-1}+1}^{e_i} P(P_k | O_{[1:k-1]}, T_i) P(R_k | P_k, O_{[1:k-1]}, T_i)$$

برآورد احتمال های شرطی رابطه بالا حتی بوسیله بزرگترین مجموعه از توالی های پروتئینی ممکن نیست به همین جهت ما  $P_k$  را فقط روی  $T_i$ ،  $P_{k-1}$  و  $P_{k-2}$  شرطی می کنیم. همچنین برای محاسبه  $P(R_k | P_k, O_{[1:k-1]}, T_i)$  از شرطی کردن روی  $O_{[1:k-1]}$  صرف نظر می کنیم. مدل های ویژه ای نیز برای حذف مشاهدات زیاد و شرطی کردن در (5) ارائه شده است که از آزمون های آماری  $\chi^2$  برای یافتن همبستگی بین مشاهدات استفاده می کند.

با محاسبه توزیع پسین  $P(m, e, T | O)$  به کمک روشهای بیزی می توانیم توزیع پسین<sup>5</sup> برای ساختار هر مشاهده،  $P(T_i | O)$ ، را بدست آوریم که  $T_i$  نشان دهنده ساختاری است که برای  $O_i$  در نظر گرفته می شود و وضعیت های  $S$ ،  $H$ ،  $C$  را اخذ می کند. برای پیش بینی ساختار  $O_i$  باید ماکزیم سازی زیر را انجام دهیم :

$$\arg \max_T P(T_i | O)$$

#### 4- نتایج :

به دلیل محدودیت های محاسباتی در اجرای شبکه های عصبی، این روش را روی مجموعه داده های RS126 تست کرده ایم. با استفاده از اطلاعات یک توالی تنها و بدون به کارگیری پروفایل انطباق چنگانه و از تقسیم توالی ها به پنج قسمت تخمین پارامترها انجام شده و نتایج بدست آمده در جدول 1 آمده است :

جدول 1: دقت پیش بینی به روش احتمال پسین

$Q_H$	$Q_E$	$Q_C$	$Q_3$
65.36	38.14	80.51	67.21

$Q_H$  و  $Q_E$ ،  $Q_C$  به ترتیب دقت پیش بینی را برای وضعیت های  $E$ ،  $C$  و  $H$  نشان می دهند و  $Q_3$  دقت پیش بینی برای هر سه وضعیت است. دقت پیش بینی بدست آمده برای  $C$  و  $Q_3$  قابل رقابت با بهترین روشهای موجود است. همچنین پیش بینی مربوط به وضعیت  $E$  که دارای همبستگی با فاصله های طولانی هستند در مقایسه با سایر روشها مناسب است.

ثابت شده است که استفاده از اطلاعات تکاملی برای توالی ها می تواند به عنوان راهی برای بهبود پیش بینی کمک بسیاری کند که در کار بعدی آن را مورد بررسی قرار می دهیم. همچنین با بزرگتر شدن مجموعه داده ها شبکه های عصبی پیش بینی اولیه بهتری را ارائه می کنند.

$O_1$ ; $O_2$ ; $O_3$	$O_4$ ; $O_5$	$O_6$ ; $O_7$	$O_8$ ; $O_9$ ; $O_{10}$ ; $O_{11}$	.....
$T_1 = C$	$e_1=3$ $T_2 = S$	$e_2=5$ $T_3 = C$	$e_3=7$ $T_4 = H$	$e_4=11$

شکل (1) : ساختار دوم پروتئین برای  $T=(C,S,C,H,...)$  و  $e=(3,5,7,11,...)$  نشان داده شده است.



قدردانی و تشکر : بخشی از هزینه این کار از محل طرح شماره  
CS1385-1-02 پژوهشگاه دانشهای بنیادی (IPM) تامین شده است .  
منابع :

- [1] Schmidler C.S, Liu J.S, and Brutlag D.L, "Bayesian Segmentation of Protein Secondary Structure," J. Computational Biology, vol. 7,nos. 1/2, pp. 233-248, 2000.  
[2] Schmidler C.S, Liu J.S, and Brutlag D.L, "Bayesian Protein Structure Prediction, " Case Studies in Bayesian Statistics, pp. 363-378, Springer, 2002.  
[3] Chu W, Ghahramani Z,Podtelezhnikov A, Wild D,"bayesian segmental models with multiple sequence alignment profile for protein secondary structure and contact map prediction ,"Proc.IEEE,Vol.3,pp. 98-113 ,2006.  
[4] Won K.J, Hamerlik T ,Bennet A.P,Krogh A." evolving hidden markov models for protein secondary structure prediction," Proc.IEEE,pp. 33-44,2005.  
[5] Aydin Z, Altunbasak Y, and Borodovsky M, "Protein Secondary Structure Prediction with Semi Markov HMMs," Proc. IEEE Int'l Conf. Acoustics Speech, and Signal Processing, 2004.

## Hidden Markov Models together with Neural Networks for Protein Secondary Structure Prediction

Hamid Pezeshk<sup>1,2,\*\*\*</sup> Mehdi Sadeghi<sup>3</sup> Changize Eslahchi<sup>4</sup> Seyed Amir Malekpour<sup>1,2</sup>

1 School of Mathematics, Statistics and Computer Science, College of Science  
and Center of Excellence in Biomathematics, University of Tehran

2 Bioinformatics Research Group, Institute for Studies in Theoretical Physics and Mathematics (IPM)

3 National Institutes of Genetic Engineering and Biotechnology

4 Faculty of Mathematical Science, Shahid Beheshti University

### Abstract:

Proteins are built from amino acids sequences. For each amino acid one of the three structures, helix (H), beta strand(S) or coil(C) are considered as protein secondary structures. For many proteins the secondary structure are not known. Protein secondary structure prediction is an important step in three dimensional structure and function determination of proteins. Neural networks are usually implemented directly for estimating emission probabilities in hidden Markov models. In this paper a new combined structure of hidden Markov models and neural networks is applied. We use the predicted structure of neural network as an initial prediction for estimating emission probabilities of amino acids in a hidden Markov model. Using this combined structure we have increased prediction precision of hidden Markov models for protein secondary structure.

**Key Words:** Hidden Markov Models, Neural Networks, Protein Secondary Structure, Bayesian Methods