

سورة

دانشگاه تهران

دانشکده ریاضی، آمار و علوم کامپیوتر

مدل های مارکفی پنهان در پیش بینی
ساختار دوم پروتئین ها

نگارش: سید امیر ملک پور نرگسی

استاد راهنما: دکتر حمید پزشک

استاد مشاور: دکتر مهدی صادقی

پایان نامه برای دریافت درجه کارشناسی ارشد
در
آمار ریاضی

شهریور ۱۳۸۷

تقدیم به

مادر و پدر مهربانم

مدل های مارکفی پنهان در پیش بینی

ساختار دوم پروتئین ها

چکیده

پروتئین ها از توالی ۲۰ اسید آمینه ساخته شده اند و برای هر یک از اسید های آمینه یکی از سه ساختار *Coil* یا *Strand*، *Helix* در نظر گرفته می شود که ساختار دوم پروتئین ها نامیده می شوند. شناسایی ساختار دوم پروتئین ها یکی از مراحل کلیدی در تعیین ساختمان سه بعدی و عملکرد پروتئین ها به شمار می رود. یافتن ساختار دوم پروتئین در درمان بیماری ها و ساخت دارو کمک بسیاری می کند. اگرچه تعیین توالی اسید های آمینه آسان است، تعیین ساختار دوم گران و وقت گیر است. به همین دلیل از روش های محاسباتی برای پیش بینی این ساختار استفاده می شود. در این تحقیق کاربرد مدل های مارکفی پنهان در پیش بینی ساختار دوم مورد مطالعه قرار می گیرد و حالت خاصی از این مدل ها که مدل نیمه مارکف قطاعی (*SSMM*) نامیده می شود مرور می شود. همچنین امکان استفاده از این مدل ها در پیش بینی *Contact Map* ها که یکی از مزیت های مدل های احتمالی است و در تعیین ساختار سه بعدی پروتئین ها کمک بسیاری می کند مورد مطالعه قرار می گیرد. این کار معمولاً به کمک روش های *MCMC* انجام می شود. از قابلیت های دیگر مدل های مارکفی پنهان توانایی این مدل ها برای در نظر گرفتن *Helical Capping Signal* ها است. همچنین امکان استفاده از این مدل ها را برای در نظر گرفتن اطلاعات تکاملی موجود در توالی های *Homolog* مرور می کنیم و از آزمون های آماری χ^2 برای یافتن همبستگی های معنی دار بین

اسیدهای آمینه استفاده خواهد شد.

مدل های *SSMM* تنها از اطلاعات و وابستگی های موجود در سمت چپ هر اسید آمینه استفاده می کنند. در اینجا روشی برای در نظر گرفتن وابستگی های چپ به راست و راست به چپ در درون توالی پروتئینی معرفی می شود. همچنین از شبکه های عصبی به منظور کاهش حجم داده ها و به کار گیری اطلاعات موجود در توالی های *Homolog* پروتئینی استفاده خواهد شد، که در آن از مدل های *SSMM* برای انجام پیش بینی نهایی استفاده می شود.

واژه های کلیدی : مدل های مارکوف پنهان، مدل های نیمه مارکوف پنهان، مدل های ممیزی، نمونه گیری به کمک زنجیر مارکوف مونت کارلو (*MCMC*)، آزمون های χ^2 ، روش های بییزی، ساختار دوم پروتئین ها.

قدردانی

با سپاس و ستایش از معبود یگانه که پرتو الطاف بی‌شمارش بر لحظه لحظه زندگی‌ام ساطع و آشکار است. تحقیق حاصل نتیجه همکاری و راهنمایی اساتید و دوستان بزرگواری است که از محبت آنها نهایت سپاس را دارم. پیش از هر چیز جای دارد قدردانی عمیق خود را نسبت به استادم دکتر حمید پزشک برای دانش، بینش، حمایت و راهنمایی‌های ایشان در طول تحصیلاتم اظهار دارم و به خاطر بینشی که از علم آمار به من داده‌اند از ایشان تشکر نمایم. دکتر پزشک در مراحل مختلف تحقیق صمیمانه همراهی‌ام کردند و در بهبود کار از هیچ تلاشی دریغ نوزیدند. بدون هدایت ایشان قادر به پیشرفت و انجام این کار نبودم.

از دکتر مهدی صادقی و دکتر چنگیز اصلاحچی که از راهنمایی‌های ارزشمند ایشان استفاده کردم و با نگاه ظریفشان زوایای مختلف تحقیق را بر من گشودند تشکر ویژه‌ای دارم. همچنین از دکتر عمید رسولیان که داوری پایان‌نامه‌ام را به عهده گرفتند سپاسگذارم.

در پایان از آقای امیر لکی زاده و خانم سیما نقی زاده کمال تشکر را دارم.

سید امیر ملک پور

شهریور ۱۳۸۷

فهرست مندرجات

| | | |
|----|---|----|
| ۱ | تاریخچه و دورنما | ۱ |
| ۱ | ۱.۱ بیوانفورماتیک ساختاری در عصر پس از ژنوم | ۱ |
| ۲ | ۲.۱ تاریخچه پیش بینی ساختار دوم پروتئین ها | ۲ |
| ۳ | ۱.۲.۱ روش های کلاسیک | ۳ |
| ۵ | ۲.۲.۱ روش های مدرن | ۵ |
| ۶ | ۳.۲.۱ روش های بیزی | ۶ |
| ۸ | ۳.۱ مطالب بررسی شده در این پایان نامه | ۸ |
| ۱۱ | ۲ مقدمه ای بر مدل های مارکفی پنهان | ۱۱ |

| | | |
|----|---|-------|
| ۱۱ | زنجیرهای مارکف | ۱.۲ |
| ۱۳ | مدل های مارکفی پنهان | ۲.۲ |
| ۱۶ | چگونه احتمال توالی مشاهده شده به شرط مدل رامحاسبه کنیم؟ | ۱.۲.۲ |
| ۱۹ | چگونه توالی وضعیت های پنهان را براساس توالی مشاهدات بیابیم؟ | ۲.۲.۲ |
| ۲۲ | تخمین پارامترهای یک <i>HMM</i> | ۳.۲.۲ |
| ۲۶ | مشکل اساسی مدل های مارکفی پنهان استاندارد | ۴.۲.۲ |
| ۲۶ | مدل های نیمه مارکف پنهان | ۳.۲ |
| ۳۲ | مقدمه ای بر پروتئین ها و نقش آنها در جریان اطلاعات ژنتیک | ۴.۲ |
| ۳۵ | ساختار اول پروتئین ها | ۱.۴.۲ |
| ۳۶ | ساختار دوم پروتئین ها | ۲.۴.۲ |
| ۳۷ | ساختار سوم پروتئین ها | ۳.۴.۲ |
| ۳۹ | ساختار چهارم پروتئین ها | ۴.۴.۲ |
| ۴۰ | لزوم پیش بینی ساختار دوم پروتئین ها | ۵.۲ |

| | | |
|----|--|--------|
| ۴۳ | مدل های نیمه مارکف پنهان درپیش بینی ساختار دوم پروتئین ها | ۳ |
| ۴۳ | مقدمه | ۱.۳ |
| ۴۵ | پیش بینی ساختار دوم با استفاده از اطلاعات تکاملی | ۲.۳ |
| ۴۶ | پروفایل انطباق چند گانه توالی ها | ۱.۲.۳ |
| ۴۸ | توزیع پیشین | ۲.۲.۳ |
| ۴۸ | تابع درستنمایی | ۳.۲.۳ |
| ۵۲ | توزیع پسین | ۴.۲.۳ |
| ۵۲ | در نظر گرفتن اثرات متقابل در فاصله های دور در β - sheet ها | ۵.۲.۳ |
| ۵۷ | β - sheet Contact Maps | ۶.۲.۳ |
| ۵۸ | تخمین پارامترها | ۷.۲.۳ |
| ۶۵ | روش نمونه گیری برای پیش بینی | ۸.۲.۳ |
| ۶۹ | معیارهای اندازه گیری دقت پیش بینی | ۹.۲.۳ |
| ۷۱ | نتایج پیش بینی | ۱۰.۲.۳ |
| ۷۷ | نتایج پیش بینی برای $contact\ map$ ها | ۱۱.۲.۳ |
| ۸۲ | پیش بینی ساختار دوم با استفاده از توالی اسید های آمینه | ۳.۳ |
| ۸۲ | آنالیز وابستگی بین اسید های آمینه | ۱.۳.۳ |
| ۸۴ | مدل نیمه مارکف قطاعی برای توالی اسید های آمینه | ۲.۳.۳ |
| ۸۶ | آموزش مدل به صورت تکراری | ۳.۳.۳ |
| ۸۷ | نتایج پیش بینی با استفاده از یک توالی تنها | ۴.۳.۳ |

۴ مدل های نیمه مارکف پنهان دوطرفه و کاربرد شبکه های عصبی در

۸۸ کاهش داده

۸۸ ۱.۴ مقدمه

۸۹ ۲.۴ مدل های دو طرفه $TS - SSMM$

۹۱ ۱.۲.۴ مدل احتمالی اول

۹۳ ۲.۲.۴ مدل احتمالی دوم

۹۳ ۳.۲.۴ پیش بینی توسط مدل $TS - SSMM$

۹۶ ۳.۴ ترکیب شبکه های عصبی و مدل های $SSMM$

۹۷ ۱.۳.۴ استفاده از شبکه های عصبی برای کاهش حجم داده ها

۹۸ ۲.۳.۴ مدل $SSMM$ به عنوان تابع تصمیم گیری

۹۹ ۳.۳.۴ مجموعه داده های آموزشی برای تخمین پارامترها و ارزیابی نتایج

۹۹ ۴.۳.۴ نتایج انجام پیش بینی ها

۱۰۷ الف پروفایل های انطباق چندگانه

۱۰۹ ب نمونه گیری از فضای $\beta - sheet$

۱۱۱ ج الگوریتم متروپلیس-هستینگ

۱۱۳ د واژه نامه ی فارسی به انگلیسی

تاریخچه و دورنما

۱.۱ بیوانفورماتیک ساختاری در عصر پس از ژنوم

نقشه ژنوم انسان در سال ۲۰۰۲ کامل شده است، علاوه بر این توالی ژنومی بسیاری از موجودات دیگر یا کد گشایی شده است و یا در حال انجام است. در دست داشتن این داده ها خبر از انقلابی در پزشکی و زیست شناسی می دهد و اساس یک بینش به سوی شناخت مکانیسم بیماری ها و گسترش روشها و داروهای جدید برای مهار آنها است و لازمه این کار تلاش بسیاری است که باید انجام شود. زیرا توالی ژنومی موجودی مانند انسان دارای ۳ میلیارد جز است که به طور تخمینی شامل ۳۵۰۰۰ ژن می شود و جایگاه عمل ژنها پروتئین ها به حساب می آیند. فرآیند کشف ساختار و عملکرد این پروتئین ها و پیشرفت های بعدی در کشف داروهای جدید وابستگی زیادی به شناسایی ساختار سه بعدی و فضایی پروتئین ها دارد. متأسفانه هرچند می توان توالی های پروتئینی را به سادگی با استفاده از توالی ژنومی تعیین کرد ولی کشف ساختار سه بعدی آنها به یکی از سخت ترین مسائلی مبدل شده است که علم تا کنون با آن روبه رو بوده است. با توجه به اینکه روش های آزمایشگاهی تعیین ساختار

بسیار گران و وقت گیر هستند گسترش روش های محاسباتی برای تعیین و پیش بینی این ساختارها بسیار ضروری به نظر می رسد. شاخه ای که به بررسی این مسائل به کمک روش های محاسباتی می پردازد بیو انفورماتیک ساختاری نام دارد. و نه تنها ساختار پروتئین ها بلکه گستره وسیعی از مسائل دیگر را نیز در بر می گیرد. با توجه به پایان یافتن پروژه ژنوم انسان، تلاش های بسیاری برای گسترش و ساخت یک مجموعه بزرگ از ساختار های دومی که با دقت بالا تعیین شده اند و انواع متفاوت ساختار های مشاهده شده در طبیعت را شامل می شوند افزایش یافته است. در راستای این هدف کشف روشهای محاسباتی و آماری جدید نیز سهم زیادی در پیشبرد آن خواهد داشت.

۲.۱ تاریخچه پیش بینی ساختار دوم پروتئین ها

پیش بینی ساختار دوم پروتئین ها یکی از مراحل کلیدی در پیش بینی ساختار فضایی و عملکرد پروتئین ها به شمار می رود. در مجموع ۲۰ نوع اسید آمینه وجود دارد و در ساختار دوم پروتئین ها به هر اسید آمینه یکی از سه ساختار *Helix*، *β -strand* و *Coil* نظیر می شود. که این ساختارها نشان دهنده نظم فضایی محلی اسید های آمینه نسبت به یکدیگر هستند. مسئله پیش بینی ساختار دوم بدون هیچ اطلاعی از ساختار فضایی انجام می شود و تصویر کردن یک فضای 2^L بعدی، که در آن L طول توالی پروتئینی است، از اسیدهای آمینه به فضای 3^L بعدی از ساختار های دوم است. از آنجایی که اطلاع در مورد ساختار دوم می تواند در الگوریتم هایی که به مسئله ساختار فضایی می پردازند به کار رود پیش بینی ساختار دوم بسیار مورد علاقه است. این مسئله در طول چند دهه توجه زیادی را به خود جلب کرده است اما هنوز به عنوان مسئله ای مشکل باقی مانده است. ما در اینجا چند نگرش کلاسیک مهم و روشهای مدرن را به طور خلاصه ذکر می کنیم.

۱.۲.۱ روش های کلاسیک

به طور ضروری همه روش های پیش بینی ساختار دوم بر اساس این قاعده هستند که اسید های آمینه مختلف دارای فراوانی های مختلفی در هر یک از ساختار های دوم هستند. این اختلاف فراوانی ها می تواند به کمک روشهای آماری مدل شود و به صورت توزیع شرطی اسید های آمینه روی ساختارهای دوم بیان شود. در میان قدیمی ترین روش های انتشار یافته برای استفاده از اختلاف فراوانی ها می توان به روش چو-فسمن^۱، [۲]، [۱] اشاره کرد. چو و فسمن برای هر اسید آمینه یک پارامتر بر اساس فراوانی مشاهده شده در هر نوع از ساختار های دوم در نظر گرفتند که بر اساس انتقال های بین *helix - strand* است اما به طور ساده به صورت زیر تعریف می شوند:

$$P_s(a) = \frac{n_{a,s}/n_{.,s}}{n_{a,.}/n_{.,.}} \quad (1.2.1)$$

که در آن a نشان دهنده اسید آمینه است و $s \in \{helix, \beta - strand, coil\}$ نشان دهنده کلاس ساختارهای دوم می باشد. $n_{x,y}$ نشان دهنده تعداد تجربی اسید آمینه x است که در ساختار از نوع y مشاهده می شود و $n_{x,.} = \sum_y n_{x,y}$ است. رابطه بالا به عنوان برآوردی از $\frac{P(a|s)}{P(a)}$ یا $P(s|a)$ با فرض یک توزیع یکنواخت برای s در نظر گرفته می شود. حاصلضرب این پارامترها برای یک توالی از اسیدهای آمینه همسایه به عنوان یک معیار ساده برای احتمال زیر توالی با فرض استقلال مشاهدات در نظر گرفته می شود. بر اساس این پارامترها ۲۰ اسید آمینه به کلاس های مختلفی دسته بندی شدند که در جدول ۱.۲.۱ مشاهده می شود. شایان توجه است که این طبقه بندی با طبقه بندی مدرن که بر اساس ویژگی های فیزیکی اسیدهای آمینه صورت می گیرد کاملاً تطابق دارد. بر اساس طبقه بندی بالا و بر اساس تعداد اسید های آمینه هر کلاس در زیر توالی های چهار تایی و پنج تایی نمره هایی به

این زیرتوالی‌ها در یک پروتئین جدید تخصیص داده شده است. ناحیه‌های با نمره‌های بالا به عنوان ناحیه‌های با پتانسیل $helix$ یا $\beta - strand$ در نظر گرفته شدند. سپس یک سری از قواعد را برای پیش‌بینی ساختار تعریف کردند. این قواعد به عنوان اولین الگوریتم برای پیش‌بینی ساختار دوم پروتئین‌ها مطرح هستند.

جدول ۱.۲.۱: تخصیص اسیدهای آمینه به کلاس‌های ساختار دوم که به وسیله چو و فسمن مورد استفاده قرار گرفته است.

| کلاس | اسید آمینه |
|-----------------------------------|----------------|
| تشکیل دهنده $helix$ قوی | <i>EAL</i> |
| تشکیل دهنده $helix$ | <i>HMQWVF</i> |
| تشکیل دهنده $helix$ ضعیف | <i>KI</i> |
| $helix$ - بی تفاوت | <i>DTSRC</i> |
| $helix$ - شکن | <i>NY</i> |
| $helix$ - شکن قوی | <i>PG</i> |
| تشکیل دهنده $\beta - strand$ قوی | <i>MVI</i> |
| تشکیل دهنده $\beta - strand$ | <i>CYFQLTW</i> |
| تشکیل دهنده $\beta - strand$ ضعیف | <i>A</i> |
| $\beta - strand$ - بی تفاوت | <i>RGD</i> |
| $\beta - strand$ - شکن | <i>KSHNP</i> |
| $\beta - strand$ - شکن قوی | <i>E</i> |

دومین روشی که برای پیش‌بینی ساختار بر اساس فروانی اسیدهای آمینه پیشنهاد شده است روش *Garnier - Osguthorpe - Robson* است که دارای نگرش تئوری اطلاعات است. و محتوای اطلاعات نهفته در هر اسید آمینه a را در ساختار دوم از نوع s محاسبه می‌کند. و آن را به صورت $I(s, a) = \log \frac{P(a|s)}{P(a)}$ می‌نویسد. این برابر با لگاریتم پارامترهایی است که توسط چو و فسمن در نظر

گرفته شدند. سپس تفاوت این اطلاعات $I(\Delta s : a) = I(s, a) - I(-s, a)$ را محاسبه کردند که به سادگی برابر با $\log \frac{P(a|s)}{P(a|-s)}$ است. همانند چو-فسمن این نسبت می تواند برای هر اسید آمینه محاسبه شود و ساختار توالی های جدید را پیش بینی کند. گرنیر و دیگران در [۳] این روش را برای پنجره های اطراف هر اسید آمینه به کار گرفتند. و برای هر اسید آمینه بر اساس مشاهدات اطرافش پیش بینی انجام دادند. و با استفاده از مستقل بودن اسیدهای آمینه در هر پنجره، تفاوت اطلاعات را در هر پنجره جمع بستند و پیش بینی برای کل توالی پروتئینی با حرکت این پنجره در طول توالی محاسبه می شود. مدل های با وابستگی های مرتبه بالا تر نیز در نسخه های جدید به کار گرفته شدند. دقت دو روش گفته شده در بالا ۵۵٪ گزارش شده اند.

۲.۲.۱ روش های مدرن

شاید یکی از مهمترین تاثیرات روش *Garnier – Osguthorpe – Robson* معرفی ساختار پنجره ای در پیش بینی باشد که مسئله را به تصویر کردن یک بردار با طول ثابت به یک کلاس از خروجی ها کاهش می دهد. گرنیر در روی کاربرد بسیاری از روشهای طبقه بندی آماری در آمار، یادگیری ماشینی، و تشخیص الگو گشود. گستره وسیعی از کاربرد این الگوها در سالهای بعد برای پیش بینی ساختار پروتئین ها به کار گرفته شده است.

بسیاری از موفق ترین این روشها بر اساس مدل های پیش بینی کننده غیر خطی مانند شبکه های عصبی و طبقه بندی بر اساس نزدیکترین همسایگی بوده اند. یکی از منفعت های این روش ها توانایی برای به دام انداختن وابستگی بین مشاهدات در همسایگی یکدیگر است. بسیاری از این روشها مانند روش *Garnier – Osguthorpe – Robson* از ساختار پنجره ای استفاده کرده اند. یکی از معروفترین روشهای پیش بینی ساختار دوم که مرز ۷۰٪ دقت را شکست، روش *PHD* است. *PHD* بر اساس

مدل شبکه های عصبی چند لایه بوده است. که دقت پیش بینی را با استفاده از انطباق چند گانه توالی های همریخت^۲ افزایش داده است. با این حال *PHD* وقتی که روی یک توالی تنها به کار می رود با دقتی در حدود ۶۳٪، دقتی بالا تر از سایر روشهای موجود ندارد. برای اطلاع بیشتر در مورد این روش مرجع [۴] را ببینید. به هر حال یکی از مزایای این روش معرفی توالی های همریخت به عنوان منبعی مناسب از اطلاعات برای پیش بینی ساختار دوم است و به طور تقریبی بهبودی در محدوده ۶ تا ۴ درصد را برای پیش بینی ساختار دوم، نسبت به روشهایی دارد که تنها از یک توالی تنها برای پیش بینی استفاده می کنند.

۳.۲.۱ روش های بیزی

استنباط بیزی یک مدل کلی را برای پیش بینی ساختار پروتئین ها با استفاده از توالی پروتئین فراهم می کند. و به صورت سودمندی از طریق گسترش ابزارهای زیر عمل می کند:

- یک مدل آماری را بر اساس رابطه های بین ساختار/ توالی فراهم می کند که پارامترهای این مدل از روی ساختارهای پروتئینی که از طریق کارهای آزمایشگاهی تعیین شده اند پیش بینی می شوند.

- روشهای محاسباتی را برای استنباط درباره ساختار پروتئین مورد استفاده قرار می دهد.

به کمک دو ابزار بالا، استنباط بیزی می تواند پیش بینی های دقیقی را برای ساختار پروتئین فراهم کند، همچنین برآورد های قابل اعتمادی را برای عدم قطعیت پیش بینی به دست می دهد. این در حالی است که برآورد عدم قطعیت پیش بینی به کمک روشهای ممیزی مانند شبکه های عصبی

ممکن نیست. به کمک مدل های بیزی و استفاده از توزیع توام ساختار-توالی می توان منابع متفاوتی از اطلاعات را در مدل لحاظ کرد. یکی دیگر از مزایای مدل های احتمالی بیزی نسبت به مدل های ممیزی در مدل کردن و پیش بینی اثرات متقابل بین β -strand ها است که به الگو در آوردن این اثرات یکی از مشکلات بزرگ در پیش بینی توپولوژی پروتئین ها به شمار می رود و مدل های ممیزی به طور کلی در مدل کردن این اثرات ناتوان یا ضعیف به شمار می روند. مدل های مارکوفی پنهان یکی از معروف ترین روشهای بیزی هستند که در زمینه های مختلفی به کار گرفته شده اند. و به عنوان ابزاری برای *sequence labeling* به شمار می روند. یکی از اولین روشهایی که برای پیش بینی ساختار دوم به کار رفته است توسط آسای^۳ و دیگران^۵، ارائه شده است که آنها چهار نوع ساختار دوم را در نظر گرفتند و هر یک از این ساختارهای دوم به عنوان زیرمدلی مطرح شدند و برای هر زیر-مدل چهار یا پنج وضعیت پنهان در نظر گرفته شد که به طور جداگانه ای با استفاده از توالی های پروتئینی از پیش طبقه بندی شده در هر یک از چهار ساختار آموزش می دیدند. سپس این زیر-مدل ها که هر یک شامل وضعیت های پنهان می شدند با یکدیگر ادغام شدند. دقت این روش ۵۴/۷٪ شده است. روش دیگری نیز از این مدل الهام گرفته و در سال ۲۰۰۶ به دقت ۶۸/۸٪ برای توالی تنها و ۷۵٪ با استفاده از توالی های همریخت رسیده است. مارتین^۴ و دیگران در [۷]، [۶] نیز برای هر ساختار از وضعیت های متفاوت و متعددی استفاده کردند و در آن از الگوریتم *EM* برای یافتن تعداد بهینه وضعیت های پنهان برای هر ساختار استفاده کردند. در دو روش گفته شده در بالا تنها یک مشاهده از هر وضعیت پنهان انتشار می یابد. یکی از دلایل استفاده از وضعیت های مختلف برای هر ساختار مدل کردن توزیع های مختلف طولی قطاع هایی است که در هر یک از ساختارهای

Asai^۳Martin^۴

پروتئینی مشاهده می‌شوند. اما اشمیدلر^۵ در [۹], [۸] به جای استفاده از وضعیت های مختلف برای هر ساختار، از مدل های مارکفی پنهان گسترش یافته استفاده کردند که به محض ورود به هر وضعیت پنهان یک توالی از مشاهدات از آنها انتشار می یابد. این مدل که به مدل نیمه مارکف قطاعی^۶ معروف شده اند تنها از توالی اسید های آمینه برای پیش بینی ساختار پروتئین استفاده می کند دارای دقتی در حدود ۶۷٪ است. مدل های بررسی شده در این پایان نامه نیز مدل های نیمه مارکف قطاعی هستند.

۳.۱ مطالب بررسی شده در این پایان نامه

در فصل دوم مدل های مارکفی پنهان مرور شده و ابزار های استنباط توسط این مدل ها شرح داده شده اند. در این فصل مدل های نیمه مارکف پنهان نیز مورد بررسی قرار گرفته اند. این فصل با مقدمه ای بر پروتئین ها و نقش آنها در جریان اطلاعات ژنتیک پایان می یابد. در فصل سوم دو مقاله مورد تجزیه و تحلیل قرار گرفته است که هر دوی آنها بر اساس مدل پیشنهادی اشمیدلر هستند. مقاله اول به دنبال بهبود در دقت پیش بینی مدل های *SSMM* با استفاده از توالی های هم ریخت است. این مقاله [۱۰] با عنوان

Chu. W., Ghahramani. Z., Podtelezhnikov. A., David. L. W. (2006). Bayesian segmental models with multiple sequence alignment profiles for protein secondary structure and contact map prediction. *IEEE/ACM transactions on computational biology and bioinformatics*, 3:98-113.

Schmidler^۵

Segmental Semi – Markov Model(SSMM)^۶

است و با استفاده از مدل های پارامتری پروفایل انطباق چند گانه پروتئین های همریخت را مدل می کند. این روش پارامتری مشکل بزرگ دیگری را نیز حل می کند. در مدل های مارکفی پنهان با افزایش طول پنجره وابستگی تعداد پارامترهای مدل به صورت نمایی افزایش می یابد که با توجه به محدودیت در داده های پروتئینی تعیین ساختار شده برآورد این پارامترها عملی نیست. مدل پارامتری در نظر گرفته شده در این مقاله به گونه ای است که با افزایش طول پنجره وابستگی تعداد پارامترهایش به طور خطی افزایش می یابد.

این مقاله همچنین به بررسی تاثیرات متقابل بین β -strand های همی پردازد که یکی از مزایای مدل های احتمالی نسبت به مدل های ممیزی مانند شبکه های عصبی است. پیش بینی این تاثیرات متقابل در بهبود پیش بینی ساختار سه بعدی پروتئین ها کاربرد زیادی دارد و دقت آنها را افزایش می دهد. همچنین مدل کردن این تاثیرات، به ویژه هنگامی که فقط از اطلاعات توالی اسید های آمینه استفاده می شود، یکی از ابزارهای ضروری در افزایش دقت پیش بینی پروتئین ها به حساب می آید. اکثر روشهایی که امروزه برای پیش بینی ساختار دوم به کار می روند این اطلاعات را نادیده می گیرند و از اطلاعات محلی استفاده می کنند که یک چهارم اطلاعات مورد نیاز را برای پیش بینی دقیق شامل می شود.

مقاله دیگری که مورد بررسی قرار گرفته ساختارهای قویتری را برای مدل کردن وابستگی بین مشاهدات در مدل های *SSMM* پیشنهاد می کند. و در آن سعی شده تا با معرفی این وابستگی هادقت پیش بینی ساختار دوم پروتئین ها با استفاده از توالی اسید های آمینه افزایش یابد. در این مقاله [۱۸] با عنوان

Aydin. Z., Altunbasak. Y., Borodovski. M. (2004). Protein secondary structure prediction with semi markov HMMs. *Proc IEEE Int'l Conf. Acoustics Speech, and Signal*

Processing .

سعی شده تا با استفاده از آزمون های آماری χ^2 همبستگی های معنی دار بین اسید های آمینه جستجو شود و از این همبستگی ها در مدل های *SSMM* استفاده شود.

در فصل چهارم دو مدل ارائه شده است. در مدل اول سعی شده تا با استفاده از اطلاعات موجود در دو طرف هر مشاهده همبستگی های بیشتری را مدل کنیم. زیرا در مدل های *SSMM* با توجه به ساختار مارکفی بودن تنها وابستگی های یک طرفه لحاظ می شود. برای استفاده از اطلاعات موجود در دو طرف هر مشاهده مدل های شرطی متفاوتی را در نظر گرفتیم که هر یک از این مدل های شرطی وابستگی بین هر مشاهده با مشاهدات اطرافش را نشان می دهند. به منظور جمع کردن تمامی این مدل های شرطی در یک مدل، برای هر یک از این مدل ها با توجه به اهمیت شان در لحاظ کردن وابستگی ها، وزن هایی در نظر گرفته شده اند.

مدل دیگری که در فصل چهارم آمده است سعی در استفاده از اطلاعات موجود در توالی های همریخت در مدل های *SSMM* دارد. برای انجام این کار از شبکه های عصبی به عنوان ابزار کاهش حجم داده ها کمک می گیرد.