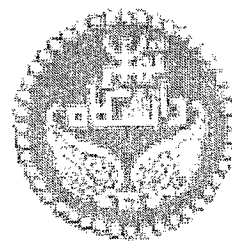
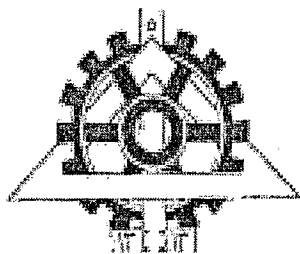


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه تهران

پردیس دانشکده‌های فنی
دانشکده مهندسی برق و کامپیوتر

تجمیع پایگاه‌های داده بیوانفورماتیک با رویکرد ترکیب داده

نگارش:

عادل اردلان

استاد راهنما: دکتر بهزاد مشیری

استاد مشاور: دکتر مسعود رهگذر

۱۳۸۷ / ۴ / ۳

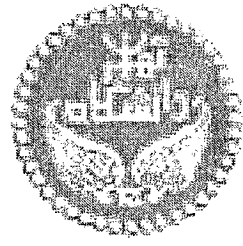
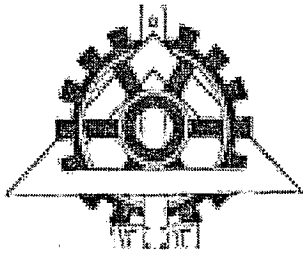
پایان‌نامه برای دریافت درجه کارشناسی ارشد

در

رشته مهندسی کامپیوتر - گرایش مهندسی فناوری اطلاعات

اسفند ماه ۱۳۸۶

۹۳۹۷۳



دانشگاه تهران

پردیس دانشکده‌های فنی

دانشکده مهندسی برق و کامپیوتر

پایان‌نامه برای دریافت درجه کارشناسی ارشد

در رشته مهندسی کامپیوتر - گرایش مهندسی فناوری اطلاعات

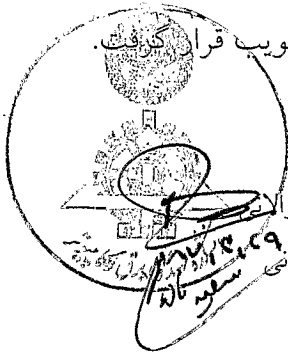
عنوان:

تجمیع پایگاه‌های داده بیوانفورماتیک با رویکرد ترکیب داده

نگارش:

عادل اردلان

این پایان‌نامه در تاریخ ۸۶/۱۱/۱۴ در مقابل هیأت داوران دفاع گردید و مورد تصویب قرار گرفت.



دکتر جواد فیض

معاونت آموزشی و تحصیلات تکمیلی پردیس دانشکده‌های فنی

دکتر پرویز جبه‌دار مارال

رئیس دانشکده مهندسی برق و کامپیوتر

دکتر سعید نادر اصفهانی

معاون پژوهشی و تحصیلات تکمیلی دانشکده مهندسی برق و کامپیوتر

دکتر بهزاد مشیری

استاد راهنما

دکتر مسعود رهگذر

استاد مشاور

دکتر مهدی صادقی

عضو هیأت داوران

دکتر بابک نجار اعرابی

عضو هیأت داوران

دکتر کارو لوکس

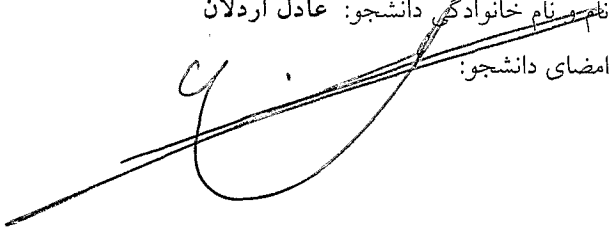
عضو هیأت داوران

تعهد نامه اصالت اثر

اینجانب عادل اردلان تأیید می‌نمایم که مطالب مندرج در این پایان‌نامه حاصل کار پژوهشی اینجانب است و به دستاوردهای پژوهشی دیگران که در این نوشته از آنها استفاده شده است مطابق مقررات ارجاع گردیده است. این پایان‌نامه قبلاً برای احراز هیچ مدرک هم سطح یا بالاتر ارائه نشده است. کلیه حقوق مادی و معنوی این اثر متعلق به پردیس دانشکده‌های فنی دانشگاه تهران می‌باشد.

نام و نام خانوادگی دانشجو: عادل اردلان

امضای دانشجو:



کامران، فہمیدہ و سحر

سہیل و حسام الدین

چکیده

تعیین ساختار مولکول‌های زیستی از مهم‌ترین و پیچیده‌ترین مسائل مطرح شده در حیطه زیست-شناسی مولکولی است. تشخیص این ساختارها به شناخت کارکرد این عوامل در مولفه‌های حیاتی کمک می‌کند. از این رو تلاش فراوانی برای بررسی قابلیت پیش‌بینی این ساختارها با استفاده از روش‌های محاسباتی صورت پذیرفته است. پیچیدگی این مساله ضرورت بررسی هرچه بیشتر را در جهت شناخت دقیق‌تر فرآیندهای حیاتی افزون می‌نماید.

در این پایان‌نامه از بستر تئوری شواهد دمپستر-شافر بعنوان مبنایی برای ارائه‌ی روشی جهت مدل‌سازی ساختار دوم پروتئین‌ها استفاده شده است. بر اساس این مدل ساختار دوم پروتئین‌ها از طریق دو منبع داده‌ای مهم مورد بررسی و ارزیابی قرار می‌گیرند و از نتایج این بررسی در پیش‌بینی ساختار پروتئین‌هایی که ساختار آنها ناشناخته است، استفاده می‌گردد. منابع اطلاعاتی مورد استفاده با توجه به این مطلب انتخاب شده‌اند که اطلاعات ساختاری قابل توجهی را در ارتباط با آرایش درشت‌مولکول‌های پروتئینی در اختیار قرار می‌دهند. این منابع عبارتند از ساختار اول پروتئین (پایگاه داده‌های پروتئین‌ها)¹ و اطلاعات جابجایی شیمیایی (پایگاه داده‌های تشدید مغناطیسی زیستی)².

در مدل ارائه‌شده رشته‌های پروتئینی بر مبنای توالی اسیدهای آمینه بررسی و شواهد موجود استخراج می‌شوند. سپس برای هر اسید آمینه اعداد جابجایی شیمیایی اندازه‌گیری شده با دقت از پیش تعیین شده، از داده‌های موجود استخراج می‌گردند. در مرحله‌ی بعد به هریک از ساختارهای مورد نظر برای انجام عملیات پیش‌بینی، جرم احتمال تخصیص داده می‌شود. نتایج بدست‌آمده نشان‌دهنده بهبود قابل ملاحظه-ای در صحت پیش‌گویی هریک از ساختارها و نیز صحت پیش‌بینی کلی (Q_3) است. از این جهت این روش می‌تواند برای پیش‌بینی قابل اطمینان ساختار دوم پروتئین‌ها مورد استفاده قرار گیرد.

کلمات کلیدی: تئوری شواهد دمپستر-شافر، بیوانفورماتیک، ساختار دوم پروتئین‌ها، طیف‌سنجی

مغناطیسی هسته.

¹ Protein Data Bank (PDB) - <http://www.rcsb.org/pdb/>

² Biological Magnetic Resonance Data Bank (BMRB) - <http://www.bmrb.wisc.edu/>

فهرست مطالب

۱۲.....	مقدمه	۱
۳.....	موضوع پایان نامه.....	۱-۱
۴.....	محدودیت های این پژوهش.....	۲-۱
۵.....	ساختار پایان نامه.....	۳-۱
۶.....	پروتئینها و ساختار دوم آنها	۲
۷.....	مقدمه	۱-۲
۱۰.....	رده بندی اطلاعات ساختاری مولکولهای زیستی.....	۲-۲
۱۰.....	ساختار اول.....	۱-۲-۲
۱۱.....	ساختار دوم.....	۲-۲-۲
۱۸.....	ساختار سوم.....	۳-۲-۲
۱۹.....	ساختار چهارم.....	۴-۲-۲
۲۰.....	پیش بینی ساختارهای پروتئین ها.....	۳-۲
۲۱.....	روش های مبتنی بر پایداری ترمودینامیکی.....	۱-۳-۲
۲۲.....	روش های مبتنی بر شیوه های آماری (مدلسازی تطبیقی).....	۲-۳-۲
۲۳.....	مطالعه چو-فسمن.....	۳-۳-۲
۲۴.....	بررسی هم ترازی دنباله ها (الگوریتم اسمیت-واترمن).....	۴-۳-۲
۲۶.....	استفاده از تکنیک های بازشناخت الگو.....	۵-۳-۲
۲۸.....	تئوری شواهد دمپستر-شافر و قوانین ترکیب	۳
۲۹.....	تئوری شواهد دمپستر-شافر.....	۱-۳
۲۹.....	چارچوب مشاهدات.....	۱-۱-۳
۲۹.....	تخصیص احتمال پایه.....	۲-۱-۳

۳۰.....	تابع باور.....	۳-۱-۳
۳۱.....	تابع محتمل بودن.....	۴-۱-۳
۳۱.....	قواعد ترکیب شواهد.....	۲-۳
۳۲.....	قاعده ترکیب شواهد دمپستر.....	۱-۲-۳
۳۴.....	قاعده ترکیب یاگر.....	۲-۲-۳

۴ طیف سنجی تشدید مغناطیسی هسته و جابجایی شیمیایی..... ۳۶

۳۷.....	طیف سنجی تشدید مغناطیسی هسته.....	۱-۴
۳۸.....	پدیده تشدید مغناطیسی هسته.....	۲-۴
۴۱.....	جابجایی شیمیایی.....	۳-۴
۴۲.....	ارتباط جابجایی شیمیایی و ساختار دوم پروتئین.....	۴-۴
۴۴.....	مدل تحلیل ترکیبی انرژی پروتئین از داده‌های جابجایی شیمیایی.....	۵-۴

۵ ارائه روش ملهم از تئوری شواهد دمپستر-شافر برای پیش بینی ساختار دوم پروتئین ها..... ۴۵

۴۶.....	مقدمه.....	۱-۵
۴۷.....	اهداف تحقیق.....	۲-۵
۴۸.....	الهام از تئوری شواهد دمپستر-شافر در پیش بینی ساختار دوم پروتئین ها.....	۳-۵
۴۹.....	آماده سازی داده ها.....	۴-۵
۴۹.....	استخراج اطلاعات فایل PDB.....	۱-۴-۵
۵۰.....	استخراج اطلاعات فایل BMR.....	۲-۴-۵
۵۰.....	آموزش.....	۵-۵
۵۰.....	استخراج شواهد.....	۱-۵-۵
۵۲.....	تشکیل چارچوب ادراک پایه.....	۲-۵-۵
۵۵.....	آزمون.....	۶-۵
۵۵.....	آزمون شواهد تک مانده ای.....	۱-۶-۵

۵۷.....	آزمون شواهد چند مانده ای	۲-۶-۵
۵۹.....	پیاده سازی ها و نتایج طبقه بندی ها	۶
۶۰.....	ساختار نرم افزار پیاده سازی شده	۱-۶
۶۱.....	بررسی حالت تک مانده ای	۲-۶
۶۲.....	استفاده از روش ساده سازی شده در امتیازدهی	۳-۶
۶۳.....	نتایج حاصل از پایگاه PDB	۱-۳-۶
۶۶.....	نتایج حاصل از پایگاه DSSP	۲-۳-۶
۶۸.....	بررسی حالت چند مانده ای	۴-۶
۶۹.....	نتایج حاصل از پایگاه PDB	۱-۴-۶
۷۲.....	نتایج حاصل از پایگاه DSSP	۲-۴-۶
۷۴.....	مقایسه و تحلیل	۵-۶
۷۷.....	نتیجه گیری و پیشنهادها	۷
۷۸.....	جمع بندی	۱-۷
۷۹.....	پیشنهاد کارهای آینده	۲-۷
۸۰.....	پیوست ۱: جابجایی های شیمیایی اندازه گیری شده برای آمینو اسید های مختلف	
۸۹.....	واژه نامه (انگلیسی - فارسی)	
۹۲.....	واژه نامه (فارسی - انگلیسی)	
۹۵.....	اختصارات	

۹۶ فهرست منابع

فهرست شکل‌ها

- شکل ۱-۲ برخی وظایف مهم پروتئین در موجودات زنده ۸
- شکل ۲-۲ ساختار سه بعدی پروتئین ۹
- شکل ۳-۲ ساختار دوم یک مولکول RNA ۱۲
- شکل ۴-۲ ساختار دوم یک مولکول پروتئین ۱۲
- شکل ۵-۲ نمای سه بعدی یک مارپیچ آلفا ۱۴
- شکل ۶-۲ نمای سه بعدی یک مارپیچ β ۱۵
- شکل ۷-۲ نمای سه بعدی یک مارپیچ پی ۱۵
- شکل ۸-۲ نمایش جاگیری رشته‌ها در صفحات بتای موازی ۱۶
- شکل ۹-۲ نمایش جاگیری رشته‌ها در صفحات بتای ضد موازی ۱۷
- شکل ۱۰-۲ نمای یک خم گاما ۱۷
- شکل ۱۱-۲ نمای یک خم بتا ۱۸
- شکل ۱۲-۲ ساختار سوم یک مولکول پروتئین ۱۹
- شکل ۱۳-۲ مقایسه ساختارهای اول تا چهارم پروتئینها ۲۰
- شکل ۱۴-۲ مثالی از هم ترازوی دو دنباله ۲۵
- شکل ۱۵-۲ نمایی از یک شبکه عصبی طبقه بندی کننده ساختار دوم پروتئین ۲۷
- شکل ۱-۳ شمایی از تعاریف تئوری شواهد ۳۱
- شکل ۱-۴ قواعد دست راست و چپ از الکترومغناطیس کلاسیک ۳۸
- شکل ۲-۴ میله مغناطیسی شبه هسته در میدان مغناطیسی ۳۹
- شکل ۳-۴ نمایی از یک طیف سنج NMR ۴۱
- شکل ۴-۴ ارتباط جابجایی شیمیایی $^1H^a$ با ساختار دوم ۴۳
- شکل ۱-۶ دیاگرام کلاس های نرم افزار پیاده سازی شده ۶۰
- شکل ۲-۶ نمونه ای از فایل پیکربندی ۶۱
- شکل ۳-۶ مقایسه روش های KNN+MV و PW+OWA ۶۴

شکل ۴-۶	صحت پیش بینی با ضریب باور خم ۱	۶۵
شکل ۵-۶	صحت پیش بینی با ضریب باور خم ۲	۶۶
شکل ۶-۶	صحت پیش بینی با ضریب باور خم ۱۰	۶۶
شکل ۷-۶	تاثیر ضریب باور خم روی صحت نتایج پیش بینی	۶۶
شکل ۸-۶	مقایسه دو روش برای ساختارهای چهارگانه در DSSP برای دقت 10^{-1}	۶۷
شکل ۹-۶	مقایسه دو روش برای ساختارهای چهارگانه در DSSP برای دقت 10^{-1}	۶۷
شکل ۱۰-۶	صحت پیش بینی با ضریب باور خم ۱	۶۸
شکل ۱۱-۶	صحت پیش بینی با ضریب باور خم ۲	۶۸
شکل ۱۲-۶	دقت پیش بینی با ضریب باور خم ۱۰	۶۸
شکل ۱۳-۶	صحت پیش بینی مارپیچ آلفا برای مقادیر مختلف S_7 برای پایگاه PDB	۶۹
شکل ۱۴-۶	صحت پیش بینی صفحه بتا برای مقادیر مختلف S_7 برای پایگاه PDB	۷۰
شکل ۱۵-۶	صحت پیش بینی مارپیچ تصادفی برای مقادیر مختلف S_7 برای پایگاه PDB	۷۰
شکل ۱۶-۶	صحت پیش بینی خم برای مقادیر مختلف S_7 برای پایگاه PDB	۷۰
شکل ۱۷-۶	صحت پیش بینی کلی برای مقادیر مختلف S_7 برای پایگاه PDB	۷۱
شکل ۱۸-۶	صحت پیش بینی مارپیچ آلفا برای مقادیر مختلف S_7 برای پایگاه DSSP	۷۲
شکل ۱۹-۶	صحت پیش بینی صفحه بتا برای مقادیر مختلف S_7 برای پایگاه DSSP	۷۲
شکل ۲۰-۶	صحت پیش بینی مارپیچ تصادفی برای مقادیر مختلف S_7 برای پایگاه DSSP	۷۳
شکل ۲۱-۶	صحت پیش بینی خم برای مقادیر مختلف S_7 برای پایگاه DSSP	۷۳
شکل ۲۲-۶	صحت پیش بینی کلی برای مقادیر مختلف S_7 برای پایگاه DSSP	۷۳

فهرست جدول‌ها

- جدول ۱-۲ نام اسیدهای آمینه و مخفف‌های آنها ۱۰
- جدول ۲-۲ مثال‌هایی از ابزارهای مبتنی بر مدل‌سازی تطبیقی ۲۳
- جدول ۱-۵ ساختارهای هشت‌گانه پروتئین‌ها ۴۸
- جدول ۲-۵ ماتریس جانشینی BLOSUM62 ۵۶
- جدول ۳-۵ شبه برنامه روش پیشنهادی ۵۷
- جدول ۱-۶ پارامترهای متغیر آزمون ۶۲
- جدول ۲-۶ پارامترهای امتیازدهی در حالت ساده سازی شده ۶۳
- جدول ۳-۶ تعداد مانده‌های موجود در هر ساختار ۶۳
- جدول ۴-۶ مقایسه نتایج طبقه‌بندی ۷۶

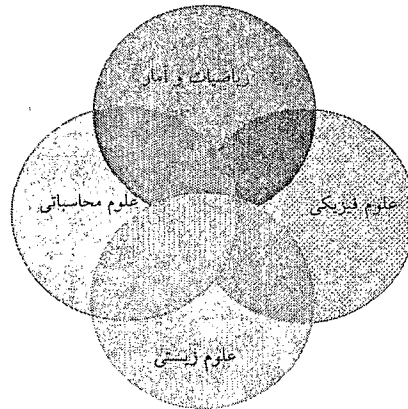
١ مقدمه

بیوانفورماتیک^۱ به ایجاد و توسعه الگوریتم‌ها، راهکارهای محاسباتی و آماری، و تئوری برای حل مسائل صوری و عملی ارائه شده و یا ملهم از مدیریت و تحلیل داده‌های زیست‌شناختی اطلاق می‌گردد. به عبارت دیگر بیوانفورماتیک علم استخراج دانش از منابع داده‌ای زیست‌شناختی به کمک تحلیل‌های کامپیوتری و حل مسایل زیست‌شناختی با استفاده از روش‌های محاسباتی و ارائه الگوریتم‌ها و روش‌های عملی و کارآ برای بررسی و نتیجه‌گیری از داده‌های زیستی است.

حجم بالای داده‌های بیوانفورماتیک از یک سو و عدم وجود مدل‌های دقیق که این سیستم‌ها را (با توجه به پیچیدگی بسیار زیاد آنها) مدل نمایند از سوی دیگر، باعث شده است که شناسایی فرآیندها و اندرکنش‌های زیست‌شناختی مولکولی جزو مسایل روز و مورد بررسی و تحقیق فراوان باشند. به عنوان مقیاسی از حجم داده‌ها در این حوزه میتوان به طول ژنوم انسان (۳،۲ میلیون جفت پایه^۲)، تعداد ژن‌های موجود (۳۲۰۰۰) و طول هر ژن (۲۷۰۰۰ جفت پایه) اشاره نمود. در این میان استفاده از روش‌های نوین آماری برای استخراج الگوهای رفتاری با استفاده از داده‌های آزمایشگاهی موجود برای شناخت و پیش-بینی فرآیندهای طبیعی با اقبال قابل ملاحظه‌ای مواجه است. همچنین افزایش روبه رشد حجم داده‌های استخراج شده در نتیجه ابداع روش‌های جدید استخراج داده‌های زیست‌شناختی منابع عظیمی از دانش را در اختیار قرار داده‌اند که استخراج این دانش روش‌های کارآمد آماری را طلب می‌کند. هدف از مطالعه بیوانفورماتیک، نه تعیین قوانین ریاضی حاکم بر فیزیک سیستم‌های زیست‌شناختی، بلکه ایجاد ابزار مورد نیاز برای تحلیل داده‌های زیست‌شناختی است. بعنوان مثال یک زیست‌شناس نیازمند ابزاری برای تشخیص ژن‌ها در *DNA* و نیز برای تخمین تفاوت میان بیان ژن‌های گوناگون در بافت‌های مختلف است. تحلیل‌هایی از این دست نیازمند ارائه‌ی مدل‌های آماری و احتمالاتی است که با توجه به طبیعت داده‌ها و نیز با تکیه بر دانش موجود بتوانند زیست‌شناس را در تشخیص یاری دهند.

¹ Bioinformatics

² Base Pair



در ایجاد داده‌های مورد استفاده در مطالعات بیوانفورماتیک سازوکارهای اتفاقی زیادی مانند فرآیندهای تصادفی دخیل در تکامل طبیعی و تصادفی بودن ذاتی فرآیندهای نمونه برداری تاثیرگذار هستند. تئوری فرآیندهای تصادفی^۱ سعی در توضیح چگونگی تغییرات فرآیندهای تصادفی در زمان یا مکان دارد. بدین ترتیب استفاده از این تئوری در مورد فرآیند تکامل طبیعی که از پیچیده‌ترین فرآیندهای تصادفی شناخته شده است، جایگاه خاصی دارد.

داده‌های بیوانفورماتیک دارای طبیعتی ناهمگون هستند؛ بدان معنا که یک مفهوم واحد در آنها در قالب‌های مختلفی بروز می‌کند. از جمله‌ی این قالب‌ها می‌توان به داده‌های رقمی، داده‌های رشته‌ای و تصاویر (دو بعدی و سه بعدی) اشاره نمود. این طبیعت ناهمگون استفاده از انبارهای داده‌های زیستی را از لحاظ برخورداری از یک دید واحد روی مفاهیم و شناسایی اندرکنش آن‌ها را با مشکل مواجه کرده است.

۱-۱ موضوع پایان نامه

در این پایان‌نامه از تئوری شواهد دمپستر-شافر بعنوان یک روش مدل‌سازی نایقینی بمنظور ارائه راه-حلی برای مسأله‌ی پیش‌بینی ساختار دوم پروتئین‌ها استفاده شده است. از این رو ابتدا مسأله‌ی مورد اشاره با استفاده از تئوری شواهد مدل‌سازی شده و ارکان اصلی مدل که عبارت از قالب شواهد، نحوه

^۱ Stochastic Process Theory

تخصیص احتمال پایه و قانون ترکیب شواهد هستند، تعیین گشته‌اند. بر مبنای این مدل شواهد موجود در داده‌های مقصد^۱ استخراج (و غنی‌سازی^۲) شده‌اند. بدین اعتبار از آنجا که داده‌ها مورد استفاده از بیش از یک منبع داده‌ای ناهمگن مورد استفاده قرار گرفته‌اند، این روش را می‌توان نوعی ترکیب داده‌ها با استفاده از روشی ملهم تئوری شواهد در نظر گرفت.

پس از تعیین بدنه‌ی شواهد مرجع^۳ از روی داده‌های ورودی آموزش، بخش آزمون بر روی شواهد استخراج‌شده از داده‌های آزمون انجام پذیرفته است. برای این منظور ابتدا سنجه‌ای برای ارزیابی میزان شباهت شواهد تعیین شده است که دارای پارامترهای قابل تغییر است. سپس با استفاده از سنجه پیش‌گفته و بر مبنای بدنه شواهد مرجع برای هر شهود آزمون، مجموعه‌ای از شواهد مرجع که بیشترین شباهت را شهود آزمون دارا هستند، تعیین گشته و نتیجه‌گیری از روی این شواهد جهت تصمیم‌سازی درباره شهود آزمون با استفاده از چند روش مرسوم ارزیابی گشته است. عمده‌ی این روش‌ها روش رای‌گیری بیشینه^۴ و میانگین‌گیری مرتب وزن‌دار^۵ بر مبنای پنجره‌ی پارزن^۶ یا نزدیکترین همسایه‌ها^۷ بوده است.

۲-۱ محدودیت‌های این پژوهش

محدودیت‌های موجود در این پژوهش عبارتند از:

- تئوری شواهد دمپستر-شافر بعنوان منبع الهام برای ارائه راه‌حل مورد استفاده قرار گرفته و گسترش‌های مطالعه‌شده این تئوری برای اعمال در این مساله مناسب تشخیص داده نشده‌اند. به عمین جهت این گسترش‌ها مورد واکاوی بیشتر قرار نگرفته‌اند.
- داده‌های مورد استفاده از پایگاه داده‌های تشدید مغناطیسی زیستی استخراج شده‌اند و در

¹ Target Data

² Enrich

³ Reference Body of Evidence

⁴ Majority Voting

⁵ Ordered Weighted Averaging (OWA)

⁶ Parzen Window

⁷ (K) Nearest Neighbors

نتیجه تمام داده‌های پروتئین‌های موجود را شامل نمی‌شوند.

۳-۱ ساختار پایان نامه

مقدمات این پایان‌نامه در فصل‌های ابتدایی مورد بررسی اجمالی قرار گرفته‌اند. در فصل دوم ساختار پروتئین‌ها به عنوان دسته مهمی از عناصر زیستی تشکیل‌دهنده ارگانیزم‌های حیاتی تبیین شده است. تئوری شواهد دمپستر-شافر و انواع قواعد ترکیب بعنوان راهبرد مورد نظر برای حل مساله مورد بحث این پایان‌نامه در فصل سوم توضیح داده شده است. طیف‌سنجی مغناطیسی هسته و جابجایی شیمیایی در فصل چهارم به اختصار مورد بررسی واقع شده‌اند. مدل پیشنهادی ملهم از تئوری شواهد در فصل پنجم واکاوی شده است و نتایج حاصل از آن در فصل ششم تشریح شده است. در انتها نتیجه‌های حاصل از این تحقیق و پیشنهادهایی برای کارهای آتی ارائه شده‌اند.

۲ پروتئین‌ها و ساختار دوم آن‌ها

۲-۱ مقدمه

در قرن هجدهم آنتوان فورکروی^۱ توانست مولکولهایی را در بدن موجودات زنده کشف کند که از لحاظ خواص شیمیایی، قدری متفاوت با دیگر مولکولها بوده و می‌توانستند دسته مجزایی را تشکیل دهند [۱]. این مولکولها برای اولین بار در موادی مانند سفیده تخم مرغ و خون یافت شدند. در سال ۱۸۳۸ به این دسته از مولکولها نام پروتئین دادند و سپس کشف کردند که این مولکولها از واحدهای کوچکتري به نام اسیدهای آمینه تشکیل شده‌اند [۱]. اما تا سال ۱۹۲۶ به نقش مهم این مولکولها در بدن یک موجود زنده پی نبرده بودند. در آن سال نشان داده شد که آنزیم *Urease* یک پروتئین است و از آنجا که وظایف مهم آنزیمها در بدن، در آن موقع شناخته شده بود، پروتئینها در رأس بسیاری از تحقیقات زیست‌شناسی قرار گرفتند [۱]. دانشمندان به سراغ استخراج کردن پروتئینها از مواد مختلف رفتند و در پی کشف ساختار فیزیکی و واکنش‌های مختلف و وظایف گوناگون این مواد، در بدن بودند. بعدها که به وسیله اشعه‌ی *X* در کریستالوگرافی، ساختار برخی پروتئینها کشف شد و ساختار پلیمری و خطی بودن این مواد و توالی دقیق واحدهای مونومری آنها (که همان اسیدهای آمینه هستند) مشخص شد، مسائل دیگری از جمله کشف ارتباط بین توالی اسیدآمینه‌ها و ساختار دو بعدی و سه بعدی پروتئین برای زیست‌شناسان مطرح شدند [۲]، [۳].

امروزه زیست‌شناسان می‌دانند که مثلاً در بدن انسان ده‌ها هزار نوع پروتئین وجود دارد [۴] که هر کدام وظایف خاصی را بر عهده دارند که برخی از این وظایف، در شکل ۱-۲ نشان داده شده‌اند. این عملکرد پروتئینها در بدن، با ساختار دو بعدی و سه بعدی و شکل فضایی آنها ارتباط دارد و این ساختار سه بعدی نیز متأثر از توالی اسیدهای آمینه‌ای است که پروتئین را می‌سازند [۵] (در شکل ۲-۲، ساختار

^۱ Antoine Fourcroy