



دانشگاه صنعتی اصفهان  
دانشکده برق و کامپیوتر

سمینار دفاع از پایان نامه کارشناسی ارشد کامپیوتر – هوش مصنوعی

پیش بینی نقشه تماس پروتئین توسط روش ماشین گروهی

ارائه دهنده

نرگس حبیبی

استاد مشاور

دکتر حسن کربکندی

استاد راهنما

دکتر محمد حسین سرایی

اردیبهشت 1388

# روند ارائه مطالب

---

- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- پیش بینی نقشه تماس
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق

# روند ارائه مطالب

---

- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- پیش بینی نقشه تماس
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق

# بیوانفورماتیک - معرفی

- تعریف

- علمی بین‌رشته‌ای (زیست‌شناسی، شیمی، فیزیک، ریاضیات و علوم کامپیوتر)
- کار بر روی داده‌های وسیع، متنوع و پیچیده زیست‌شناسی
- با هدف حل مسائل زیست‌شناسی در سطح مولکولی

- دو زمینه اصلی

- ژنومیک: تجزیه و تحلیل ژن‌های موجودات زنده
- پروتئومیک: بررسی پروتئین‌های موجودات زنده

# بیوانفورماتیک-اهداف

---

## • اهداف

- سازمان دهی داده ها
  - مثال: بانک داده پروتئین
- توسعه ابزارهایی برای کمک به تحلیل داده ها
  - مثال: مقایسه یک پروتئین با دیگر پروتئین ها
- کشف الگوها و تفسیر نتایج
  - مثال: پیش بینی ساختار پروتئین

# بیوانفورماتیک-کاربردها

## • کاربردها

### — در زیست شناسی

- یافتن ژن ها در میان ژنوم
- پیش گوئی ساختار و عملکرد محصولات ژنی
- یافتن ارتباط میان ژن ها و توالی های پروتئین

### — در علوم پزشکی

- طراحی منطقی داروها
- شناسایی علل ژنتیکی بیماری ها

# روند ارائه مطالب

---

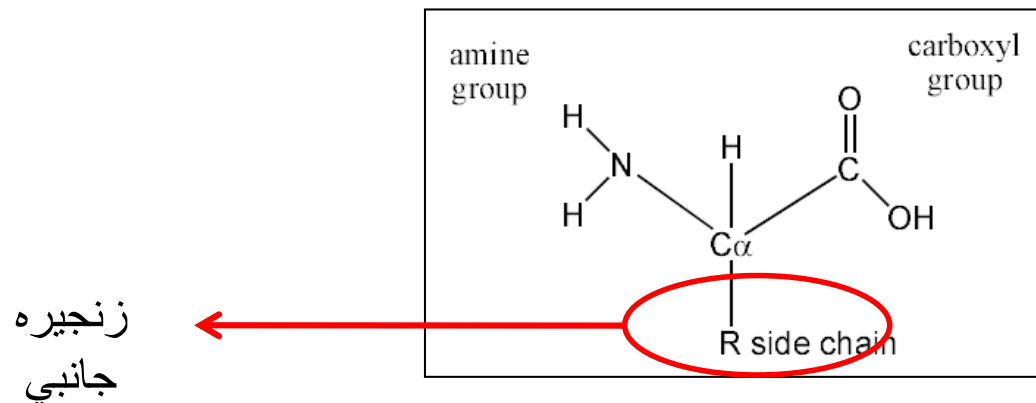
- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- پیش بینی نقشه تماس
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق

- تعریف: مولکولی مرکب، متشکل از زنجیره ای از اسیدهای آمینه
- ضرورت: مشارکت در هر پروسه درون سلولی
- نقش ها
  - آنزیمی: کاتالیز واکنش های مورد نیاز موجود زنده
  - ساختاری
    - دیواره سلول
    - غشای سلول
    - سیتوپلاسم
    - ...



## پروتئین (ادامه)

- اسیدهای آمینه
  - انواع: 20 نوع
  - ساختار

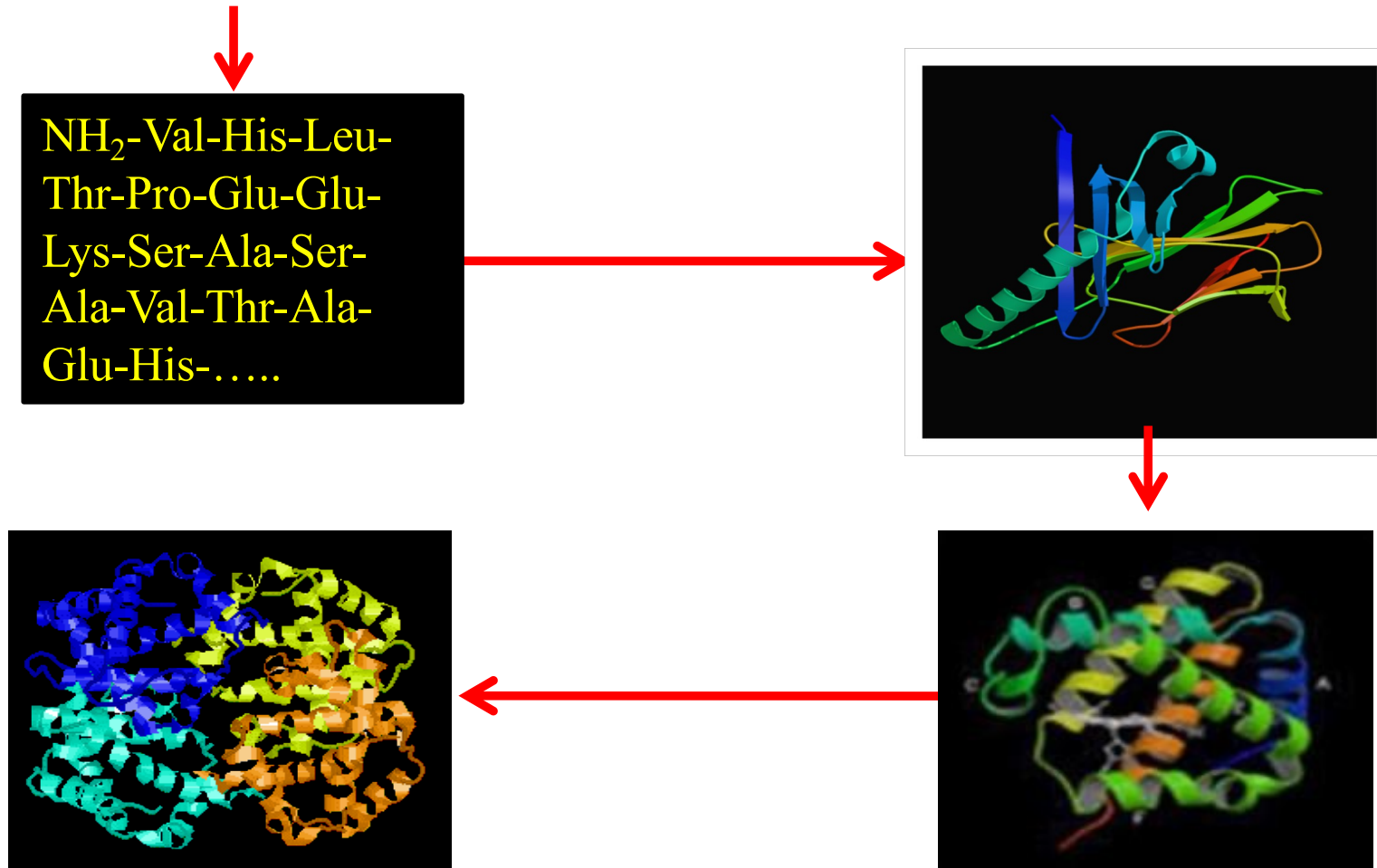


ساختار عمومی یک اسید آمینه

## پروتئین (ادامه)

- سطوح مختلف ساختار پروتئین
  - ساختار اول: توالی اسیدهای آمینه
  - ساختار دوم: مارپیچ های آلفا، صفحات بتا و دیگر ساختارهای دوم
  - ساختار سوم: حاصل تا شدن زنجیره اسیدهای آمینه، شامل چندین ساختار دوم
  - ساختار چهارم: اجتماع چند مولکول پروتئین
- نکته: وابستگی ساختارهای دوم، سوم و چهارم پروتئین به ساختار اول آن

## پروتئین (ادامه)



سطوح مختلف ساختار  
پروتئین

# روند ارائه مطالب

---

- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- پیش بینی نقشه تماس
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق

# مفاهيم محاسباتي پيش زمينه-MSA

## Multiple Sequence Alignment :MSA •

Q5E940_BOVIN	-----MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQIIRMSLRGK-AVVLVGKNTMMRKAIRGHLENN--PALE	76
RLA0_HUMAN	-----MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQIIRMSLRGK-AVVLVGKNTMMRKAIRGHLENN--PALE	76
RLA0_MOUSE	-----MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQIIRMSLRGK-AVVLVGKNTMMRKAIRGHLENN--PALE	76
RLA0_RAT	-----MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQIIRMSLRGK-AVVLVGKNTMMRKAIRGHLENN--PALE	76
RLA0_CHICK	-----MPREDRATWKSNYFMKIIQLDDYPKCFVVGADNVGSKOMQIIRMSLRGK-AVVLVGKNTMMRKAIRGHLENN--PALE	76
RLA0_RANSY	-----MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQIIRMSLRGK-AVVLVGKNTMMRKAIRGHLENN--SALE	76
Q7ZUG3_BRARE	-----MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQIIRMSLRGK-AVVLVGKNTMMRKAIRGHLENN--PALE	76
RLA0_ICTPU	-----MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKOMQIIRMSLRGK-AVVLVGKNTMMRKAIRGHLENN--PALE	76
RLA0_DROME	-----MVRENKAAWKAQYFIKVVLFDFEPKCFIVGADNVGSKOMQIIRMSLRGK-AVVLVGKNTMMRKAIRGHLENN--PQLE	76
RLA0_DICDI	-----MSGAG-SKRKKLFIEKATKLFITYDKMIVAEADFGSSQLQKIRKSIRGI-GAVLMGKNTMIRKVIRDLADSK--PELD	75
Q54LP0_DICDI	-----MSGAG-SKRKNVFIEKATKLFITYDKMIVAEADFGSSQLQKIRKSIRGI-GAVLMGKNTMIRKVIRDLADSK--PELD	75
RLA0_PLAF8	-----MAKLSKQKKQMYKREKPIPEWKTLMLELELFSSKHVVFLFADLTGTPTFVVRVRKKLWKK-YPMVMVAKKRIILRAMKAAGLE---LDDN	76
RLA0_SULAC	-----MIGLAVTTTKKIAKWKVDEVAELTEKLKTHKTIITANIEGFPADKLHEIRKKLRGK-ADIKVTKNLNFNIALKNAG----YDTK	79
RLA0_SULTO	-----MRIMAVITQERKIAKWKIEEVKELEKLREYHTIIITANIEGFPADKLHDIRKKMRGM-AEIKVTKNLTFGIAAKNAG----LDVS	80
RLA0_SULSO	-----MKRLALALKQKQKVASWKEEVKELTELIKNSNTILIGNLEGFPADKLHEIRKKLRGK-ATIKVTKNLTFKIAAKNAG----IDIE	80
RLA0_AERPE	MSVSVLVGQMYKREKPIPEWKTLMLELELFSSKHVVFLFADLTGTPTFVVRVRKKLWKK-YPMVMVAKKRIILRAMKAAGLE---LDDN	86
RLA0_PYRAE	MMLAIGKRRYVTRQYYPARKVKIVSEATELLQKYPPVFLFDLHGLSSRIHEHYRRLRY-GVIKIIKPTLFKIAFTKVYGG---IPAE	85
RLA0_METAC	-----MAEERHHTHEIPQWKDEIENIKELIQSHKVFQMVGIEGILATKMKIRRDLDKV-AVLKVSNTLTERALNQLG----ETIP	78
RLA0_METMA	-----MAEERHHTHEIPQWKDEIENIKELIQSHKVFQMVGIEGILATKMKIRRDLDKV-AVLKVSNTLTERALNQLG----ESIP	78
RLA0_ARCFU	-----MAAVRGS---PPEYKVRAVEEIKRMISSEPVVAIVSFERNVPAQGMQKIRREFRGK-AEIKVVKNLTERALDALG----GDYL	75
RLA0_METKA	MAVKAKGQPPSGYEPAKVAEWRREVKELELMDEYENVGLVDLEGIPAPQLQEIIRAKLRERDIIIRMSRNTLMRIALEEKLDER--PELE	88
RLA0_METTH	-----MAHVAEWKKKEVQELHDLIKGYEVVGIANLADIPARQLQKMRQTLRDS-ALIRMSKKTLLISLAEKAGREL--ENVV	74
RLA0_METTL	-----MITAESEHKIAPWKIEEVNKLKELKNGQIIVALVDMMEVPAPQLQEIIRDKIR-GTMTLKMSRNTLIERAIKEVAEETGNPEFA	82
RLA0_METVA	-----MIDAKSEHKIAPWKIEEVNALKELKLSANVIALDMMEVPAPQLQEIIRDKIR-DQMTLKMSRNTLIKRAVEEVAEETGNPEFA	82
RLA0_METJA	-----METKVKAHVAPWKIEEVKTLKGLIKSKPVVAIVDMMDVPAPQLQEIIRDKIR-DKVKLMSRNTLIIRALKEAAEELNNPKLA	81
RLA0_PYRAB	-----MAHVAEWKKKEVEELANLKSYPVIALVDVSSMPAYPLSQMRRLLIRENGLLRVSNTLIELAIKKAAGELGKPELE	77
RLA0_PYRHO	-----MAHVAEWKKKEVEELAKLKSYPVIALVDVSSMPAYPLSQMRRLLIRENGLLRVSNTLIELAIKKAAGELGKPELE	77
RLA0_PYRFU	-----MAHVAEWKKKEVEELANLKSYPVIALVDVSSMPAYPLSQMRRLLIRENGLLRVSNTLIELAIKKAAGELGKPELE	77
RLA0_PYRKO	-----MAHVAEWKKKEVEELANLKSYPVIALVDVAGVPAYPLSKMRDKLR-GKALLRVSNTLIELAIKKAAGELGQPELE	76
RLA0_HALMA	MSAESERKTETIPEWKQEEVDVAIVMIESYESVGVVNIAGIPSRQLDMRRDLHGT-AELRVSNTLIERALDDVD---DGLE	79
RLA0_HALVO	MSESEVRQTEVIPQWKREEDLVDFIESYESVGVVAGIPSRQLDMRRDLHGS-AAVRMSRNTLVNRRALDEVN---DGFE	79
RLA0_HALSA	MSAEEQRTTEEVPEWKQEEVAELVDLLETYSVGVVNTGIPSKQLDMRRGLHGO-AALRMSRNTLLVRALEEAG---DGLD	79
RLA0_THEAC	-----MKEVSQKKELVNEITRIKASRSVAIVDTAGIRTRQIDIRGKNRGK-INLKVIKKTLLFKALENLGD---EKLS	72
RLA0_THEVO	-----MRKINPKKKEIVSELAADITKSKAVIVDIKGVTRMODIRAKNRDK-VKIKVVKKTLFKALDSIND---EKLT	72
RLA0_PICTO	-----MTEPAQWKIDFVKNLNENSRKVAIVSIKGLRNNFQKIRNSIRDK-ARIKVSRRARLLRLAIENTGK---NNIV	72
ruler	1.....10.....20.....30.....40.....50.....60.....70.....80.....90	

بخشي از يك MSA براي چند توالي

در و تئذ

## مفاهيم محاسباتي پيش زمينه-اثر فراموشي خطرناک

---

- تعريف: از بين رفتن اطلاعات قبلي يادگرفته شده توسط شبکه عصبي، با يادگيري اطلاعات جديد
- علت: مناسب نبودن وزن هاي جديد شبکه براي تشخيص الگوهاي قبلي
- يك راه حل: تکرار

# مفاهيم محاسباتي پيش زمينه-ماشين گروهی

- تعريف

- گروهی از یادگیرها
- آموزش یادگیرها بر روی یک مجموعه داده آموزشی
- ترکیب خروجی‌های اعضا

- ایده: قاعده ((تقسیم کن و پیروز شو))

- مزیت‌ها:

- بهتر بودن کارایی گروه از کارایی هر یک از اعضا
- کاهش هزینه محاسباتی با تقسیم فضای ورودی

# مفاهيم محاسباتي پيش زمينه-ماشين گروهی

## (ادامه)

- علت استفاده: ضعف مدل های یادگیری ماشین در برخورد با مقادیر زیاد داده

– رشد غیرخطی هزینه محاسباتی

– اثر نامطلوب بر قابلیت تعمیم مدل

- محدود بودن قابلیت یک مدل
- نیاز به افزایش پارامترهای مدل با افزایش تنوع الگوها
- اثر نامطلوب افزایش پارامترها بر قابلیت تعمیم مدل و بر هزینه محاسباتی



# مفاهيم محاسباتي پيش زمينه-ماشين گروهی

## (ادامه)

---

### • انواع

#### – ایستا

- میانگین گیری گروه
- بوستینگ

#### – پویا

- ترکیب یادگیرها
- ترکیب سلسله مراتبی یادگیرها

# مفاهيم محاسباتي پيش زمينه-ماشين گروهی

## (ادامه)

- مثالي از روش میانگین گیری گروه
  - گروهی از شبکه های عصبی با وزن های اولیه مختلف
  - آموزش شبکه ها با یک مجموعه داده
  - میانگین گیری از پاسخ اعضا
- مزیت های روش میانگین گیری گروه
  - افزایش کارایی
  - کاهش زمان آموزش
  - کاهش خطر بیش پوشش

# روند ارائه مطالب

---

- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- پیش بینی نقشه تماس
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق

## نقشه تماس پروتئین

- تعریف: ماتریسی دو بعدي، نشان دهنده وجود تماس یا عدم تماس بین هر دو اسید آمینه در فضا
- مفهوم تماس: کمتر بودن فاصله بین دو اسید آمینه در فضا از یک مقدار آستانه
- بیان ریاضی:

$$C = \begin{cases} 1 & D_{ij} < thr \\ 0 & otherwise \end{cases} \quad i, j = 1, \dots, N$$

## نقشه تماس پروتئین (ادامه)

- تعریف فاصله بین دو اسید آمینه

- مینیمم فاصله بین تمام اتم‌ها

- فاصله بین اتم‌های  $C_\alpha$

- فاصله بین اتم‌های  $C_\beta$

- ...

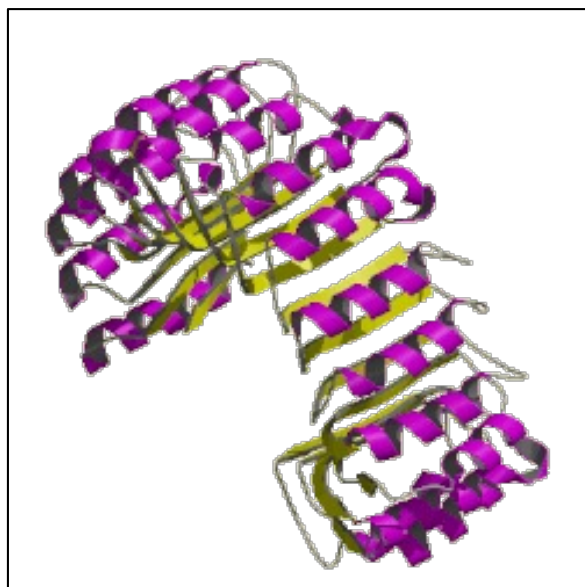
- ویژگی های نقشه تماس

- محلی بودن بیشتر تماس ها

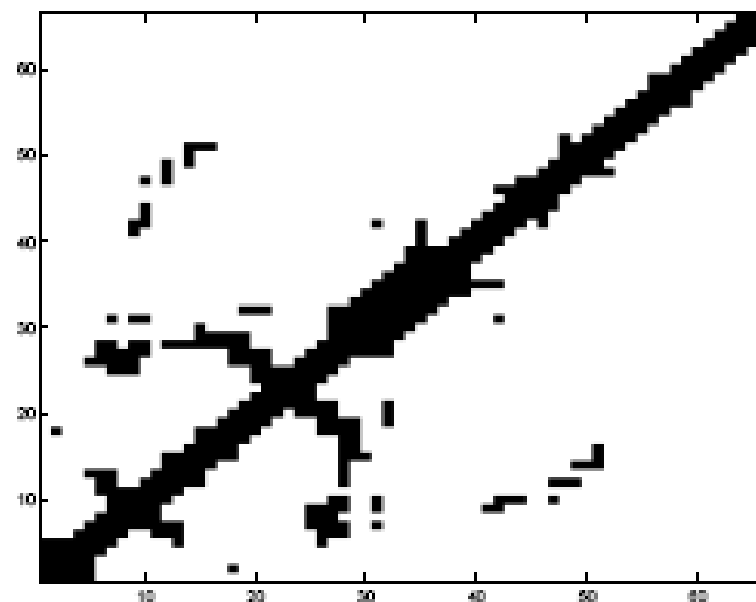
- نسبت تعداد جفت های ((تماس دار)) به ((بدون تماس)): 1 به

20

# نقشه تماس پروتئین (ادامه)



ساختار سوم یک پروتئین



نقشه تماس متناظر

# روند ارائه مطالب

---

- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- **پیش بینی نقشه تماس**
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق

## پیش بینی نقشه تماس

- هدف: پیش بینی نقشه تماس یک پروتئین بر اساس ساختار اول آن

- اهمیت

- وابسته بودن عمل پروتئین به ساختار سوم آن
- هزینه بر بودن و زمان بر بودن روش های آزمایشگاهی تعیین ساختار سوم
- نیاز به روش هایی برای تعیین ساختار سوم، صرفاً بر اساس ساختار اول
- امکان پیش بینی ساختار سوم با داشتن نقشه تماس پیش بینی شده



## پیش بینی نقشه تماس (ادامه)

- معیارهای ارزیابی

- دقت: درصد تماس های پیش بینی شده صحیح (TP)
- پوشش: درصد تماس های واقعی پیش بینی شده
- بهبود نسبت به پیش بینی کننده تصادفی
- ...

- روش های ارائه شده

- استفاده از مدل های یادگیری مختلف
- استفاده از ویژگی های گوناگون
- استفاده از معیارهای متفاوت برای ارزیابی ~~من~~ مشکل بودن مقایسه روش ها

# روند ارائه مطالب

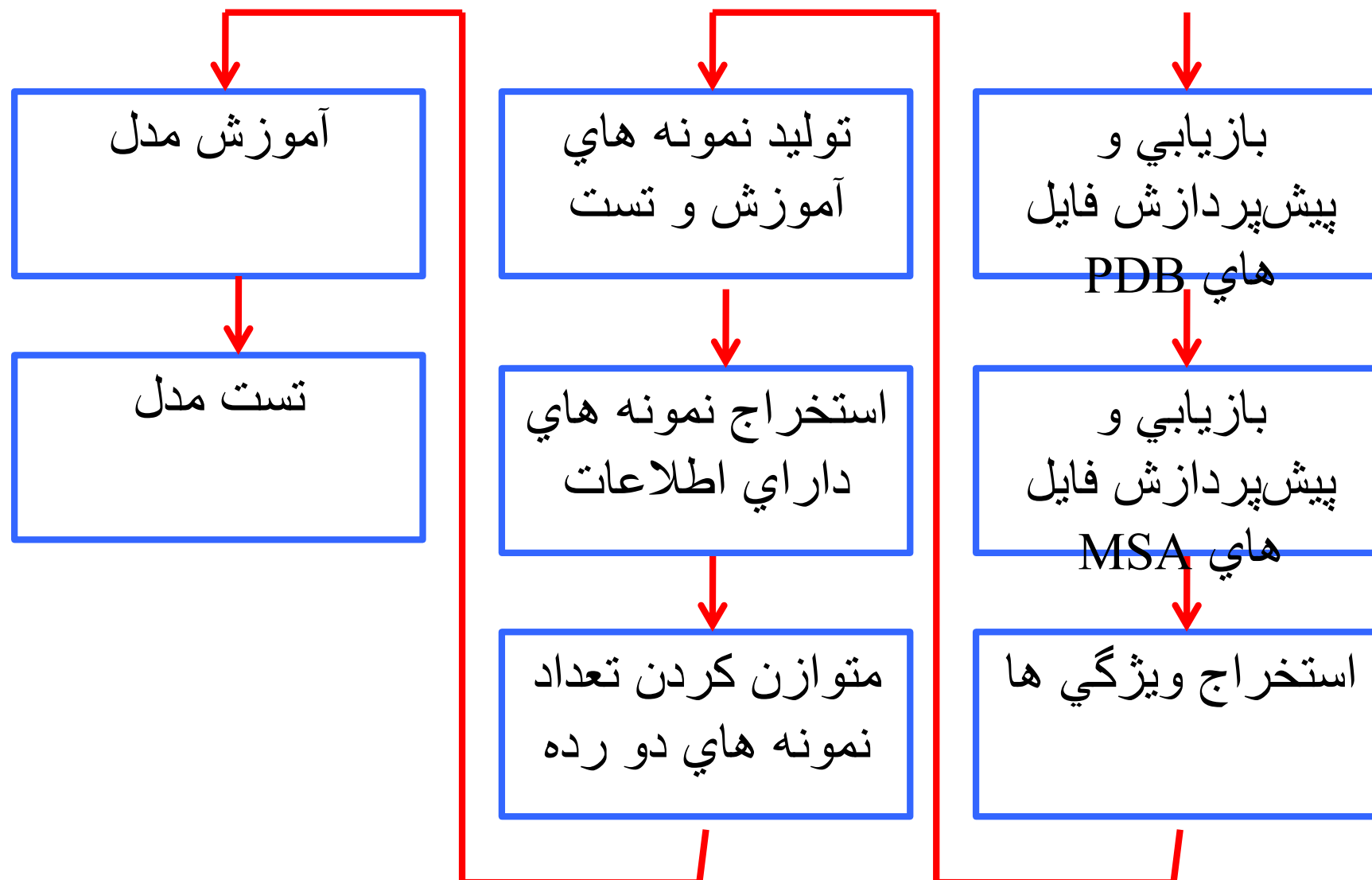
---

- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- پیش بینی نقشه تماس
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق

# روش پیشنهادی-معرفی کلی

- براساس روش ماشین گروهی
- ارائه یک ماشین گروهی جدید
- انتخاب شبکه های عصبی به عنوان اعضای گروه
- ترکیب پاسخ اعضا به صورت محاسبه میانگین پاسخ ها
- طراحی شبکه و ویژگی ها براساس روش Karplus و Shackelford
- در نظر گرفتن مساله پیش بینی نقشه تماس به عنوان یک مساله رده بندی

## روش پیشنهادی-معرفی کلی (ادامه)



## روش پیشنهادی-داده ها

- بازیابی فایل های داده از پایگاه داده PDB
- انتخاب فایل ها بر اساس لیست pdbselect-25%
- پیش پردازش فایل های PDB و استخراج مختصات اتم  $C_{\beta}$  های هر اسید آمینه
- در نظر گرفتن 80% از فایل ها برای آموزش و 20% برای تست

## روش پیشنهادي-MSA

---

- استخراج اکثر ویژگی ها از MSA پروتئین ها
- نحوه به دست آوردن: وب سرویس SAM-T08
- پیش پردازش فایل هاي به دست آمده

# روش پیشنهادی-ویژگی های مورد استفاده

---

- ویژگی های توالی
  - طول توالی - فاصله توالی
- ویژگی های تک ستونی
  - توزیع اسیدهای آمینه - آنترופی - ساختار دوم پیش بینی شده
- ویژگی های جفت ستونی
  - تعداد جفت ها - آنترופی مشترک - تمایل - MI E-value

# روش پیشنهادی-آماده سازی داده ها

---

- شامل سه مرحله
  - تولید نمونه ها
  - استخراج نمونه های دارای اطلاعات
  - متوازن کردن نسبت نمونه های دو رده



# روش پیشنهادي-آماده سازي داده ها (ادامه)

## • توليد نمونه ها

- در نظر گرفتن فاصله ميان اتم هاي
- آستانه تماس: 8 انگستروم
- آستانه تماس هاي محلي: 8
- استفاده از پنجره اي به عرض 2 براي ويژگي هاي تک ستوني
- استفاده از رتبه مقدار به جاي خود مقدار براي ويژگي هاي جفت ستوني
- اندازه بردار ورودي شبکه هاي عصبي: 419

# روش پیشنهادی-آماده سازی داده ها (ادامه)

- بردار ورودی
  - لگاریتم طول توالی
  - لگاریتم فاصله توالی
  - توزیع اسید آمینه  $i \pm 2$  (هر مکان 20 مقدار)
  - توزیع اسید آمینه  $j \pm 2$  (هر مکان 20 مقدار)
  - ساختار دوم پیش بینی شده  $i \pm 2$  (هر مکان 21 مقدار)
  - ساختار دوم پیش بینی شده  $j \pm 2$  (هر مکان 21 مقدار)
  - آنترופی  $i$
  - آنترופی  $j$
  - تعداد جفت های  $(i,j)$
  - مقدار MI E-value  $(i,j)$
  - لگاریتم رتبه MI E-value  $(i,j)$
  - لگاریتم رتبه تمایل  $(i,j)$
  - لگاریتم رتبه آنترופی مشترک  $(i,j)$
  - رده  $(i,j)$  (خروجی مطلوب)

# روش پیشنهادی-آماده سازی داده ها (ادامه)

- استخراج نمونه های دارای اطلاعات

- علت انجام

- کاهش هزینه محاسباتی
- کاهش مشکل عدم توازن تعداد نمونه های مثبت و منفی

- توجه: حذف نمونه های منفی بدون اثر بر روی مرز جداسازی دو رده

- روش

- مرتب سازی نمونه ها بر اساس چند ویژگی جفت ستونی
- انتخاب  $4L$  بالاترین نمونه ها از هر مجموعه ( $L$  طول پروتئین)
- اجتماع نمونه های انتخاب شده از هر مجموعه

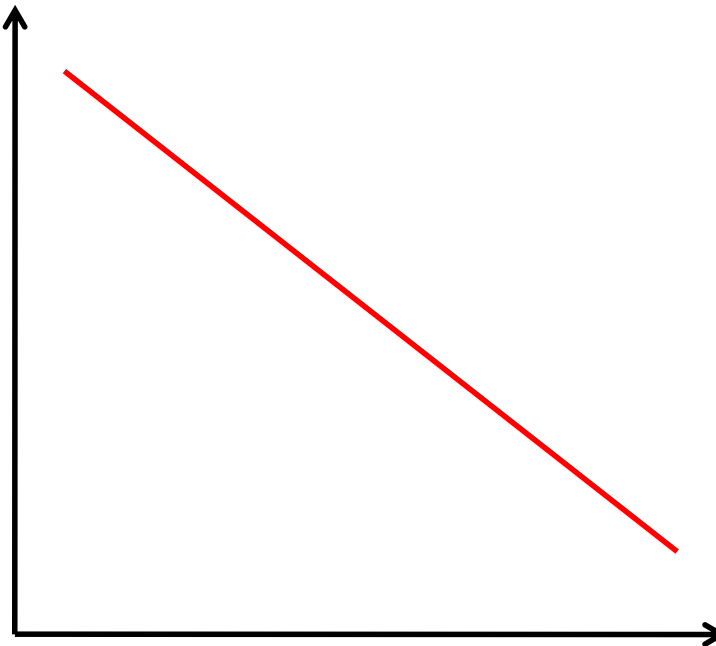
## روش پیشنهادی-آماده سازی داده ها (ادامه)

- متوازن کردن نسبت تعداد نمونه های مثبت و منفی
  - علت: بایاس شدن شبکه های عصبی به سمت رده منفی
  - نسبت متوازن سازی: 1 (مثبت) به 4 (منفی)
  - مشکل: دیده نشدن مقدار زیادی از الگوهای رده منفی در حین آموزش
  - نکته: عدم متوازن نمودن داده های تست

# روش پیشنهادی-نحوه ارزیابی کارایی

- معیار: نسبت تماس‌های پیش‌بینی شده صحیح به تعداد کل پیش‌بینی‌ها  $TP/(TP+FP)$ ، در مقابل تعداد کل پیش‌بینی‌ها  $(TP+FP)$

$$TP/(TP+FP)$$



$$\text{Log}((TP+FP)/\text{پروتئین})$$

## روش پیشنهادي-نحوه ارزیابي کارائي (ادامه)

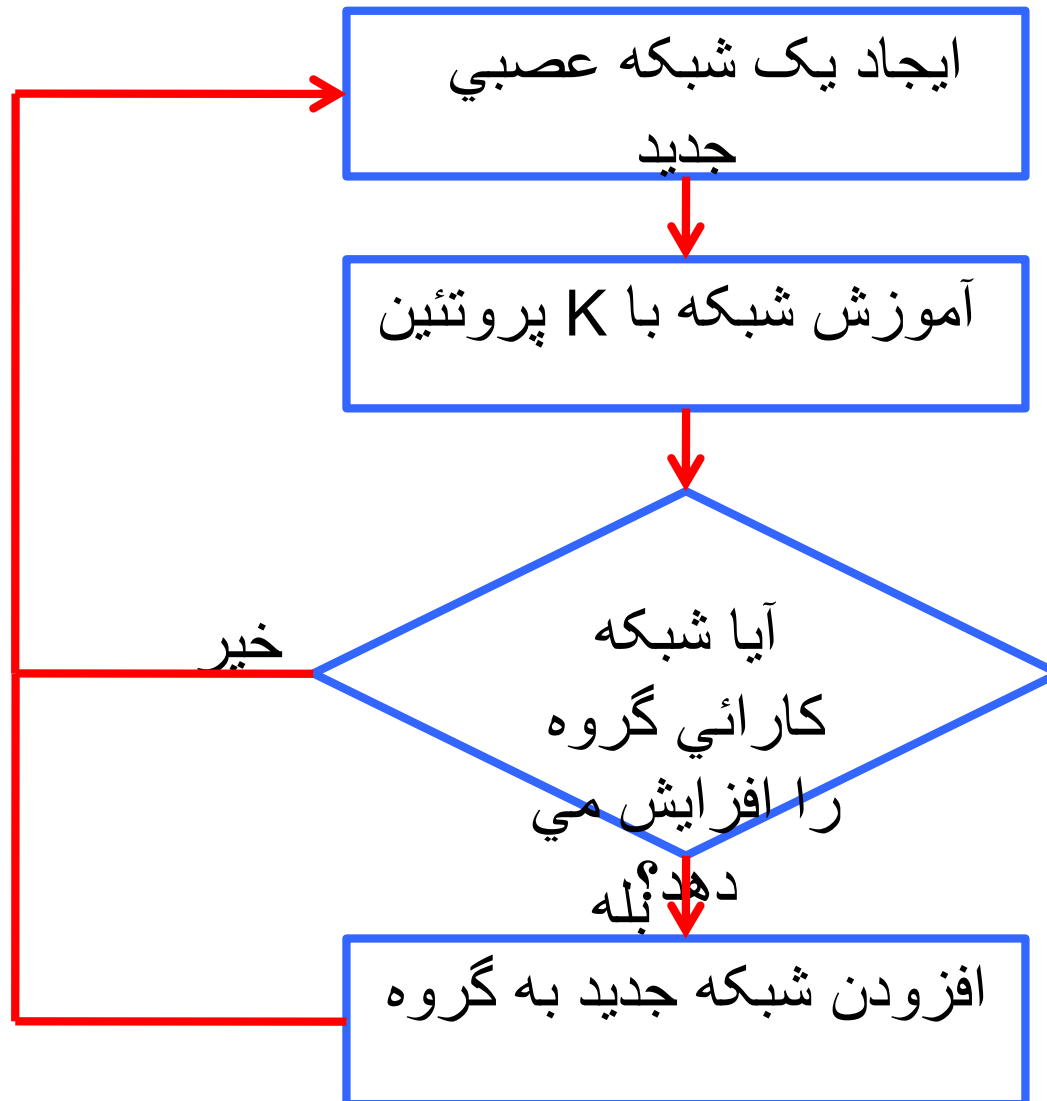
- دو رویکرد موجود براي تعیین رده هر جفت اسید آمینه
  - آستانه گذاري خروجي شبکه عصبي
  - ارزیابي کارائي برحسب آستانه هاي مختلف
- استفاده از رویکرد دوم
- نکته: نیاز کاربردهاي مختلف به تعداد پیش بيني هاي صحيح متفاوت

# روش پیشنهادی-آموزش مدل

---

- ایجاد یک ماشین گروهی ایستا در دو مرحله
  - مرحله اول: ایجاد گروه 1
  - قرار دادن آخرین عضو گروه 1 به عنوان اولین عضو گروه 2
  - مرحله دوم: ایجاد گروه 2 (گروه اصلی)

# روش پیشنهادي-آموزش مدل (ادامه)





# روش پیشنهادي-آموزش مدل (ادامه)

- ساختار شبکه هاي عصبی
  - feed-forward
  - یک لایه پنهان
  - تعداد نوروں ها: 1-45-419
  - الگوریتم آموزش: بازگشت به عقب انعطاف پذیر
  - تابع کارائی: MSE
  - حداکثر تعداد epoch هاي آموزش: 500

## روش پیشنهادي-آموزش مدل (ادامه)

- آموزش هر شبکه با  $K$  پروتئين تصادفي (تقسيم فضاي ورودي)  
– چگونگي انتخاب مقدار  $K$

- عدم آموزش يك شبکه با  $K$  پروتئين به طور سريال براي اجتناب از مشكل اثر فراموشي خطرناك

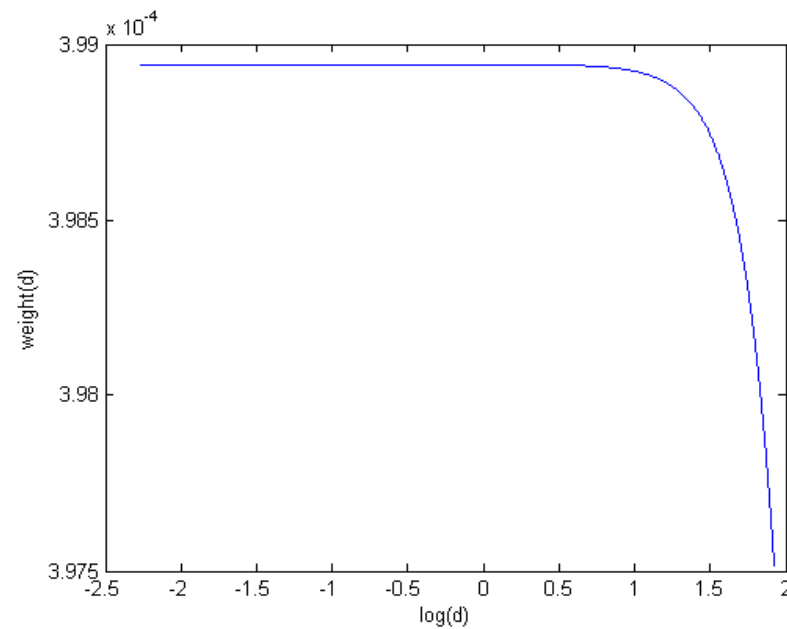
- استفاده از 90% داده هاي آموزشي براي آموزش و 10% براي جلوگيري از بيش پوشش

## روش پیشنهادي-آموزش مدل (ادامه)

- نحوه تصميم گيري درباره اضافه شدن شبکه جديد به گروه
  - محاسبه دو نمودار کارائي
    - نمودار 1: گروه جاري
    - نمودار 2: گروه جاري + شبکه جديد
  - محاسبه نمودار تفاضل نمودار 1 و نمودار 2
  - محاسبه مجموع وزن دار نمودار تفاضل
  - اضافه کردن شبکه جديد به گروه در صورت مثبت بودن مجموع

# روش پیشنهادي-آموزش مدل (ادامه)

- تابع وزن پيشنهادي



تابع وزن پيشنهادي براي محاسبه  
مجموع وزن دار نمودار تفاضل

## روش پیشنهادی-آموزش مدل (ادامه)

- نکته: تقسیم مجموعه داده آموزشی به دو زیرمجموعه آموزش و تست

- استفاده از زیرمجموعه آموزش برای آموزش هر شبکه جدید
- استفاده از زیرمجموعه تست برای تصمیم گیری در مورد اضافه شدن شبکه جدید به گروه
- عدم استفاده از زیرمجموعه تست برای تست نهایی مدل

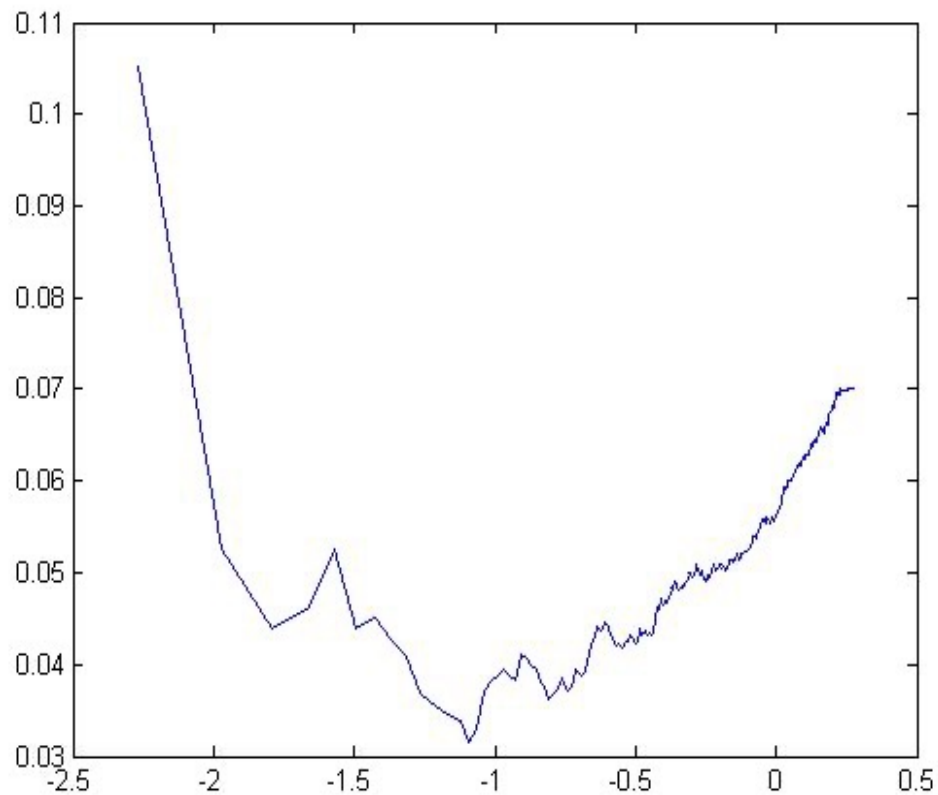
## روش پیشنهادی-تست مدل

- کارایی مدل: میانگین گیری از کارایی های مدل بر روی هر پروتئین تست
- پیاده سازی سه روش پیش بینی
  - روش پایه: آموزش یک شبکه با تمام داده های آموزشی به طور سریال
  - روش میانگین گیری گروه: آموزش 10 شبکه (با وزن های اولیه متفاوت) با تمام داده های آموزشی به طور سریال و میانگین گیری از پاسخ شبکه ها
  - روش پیشنهادی: گروهی با 105 عضو

# روش پیشنهادي-بررسي نتايج

- نتيجہ روش پایه:

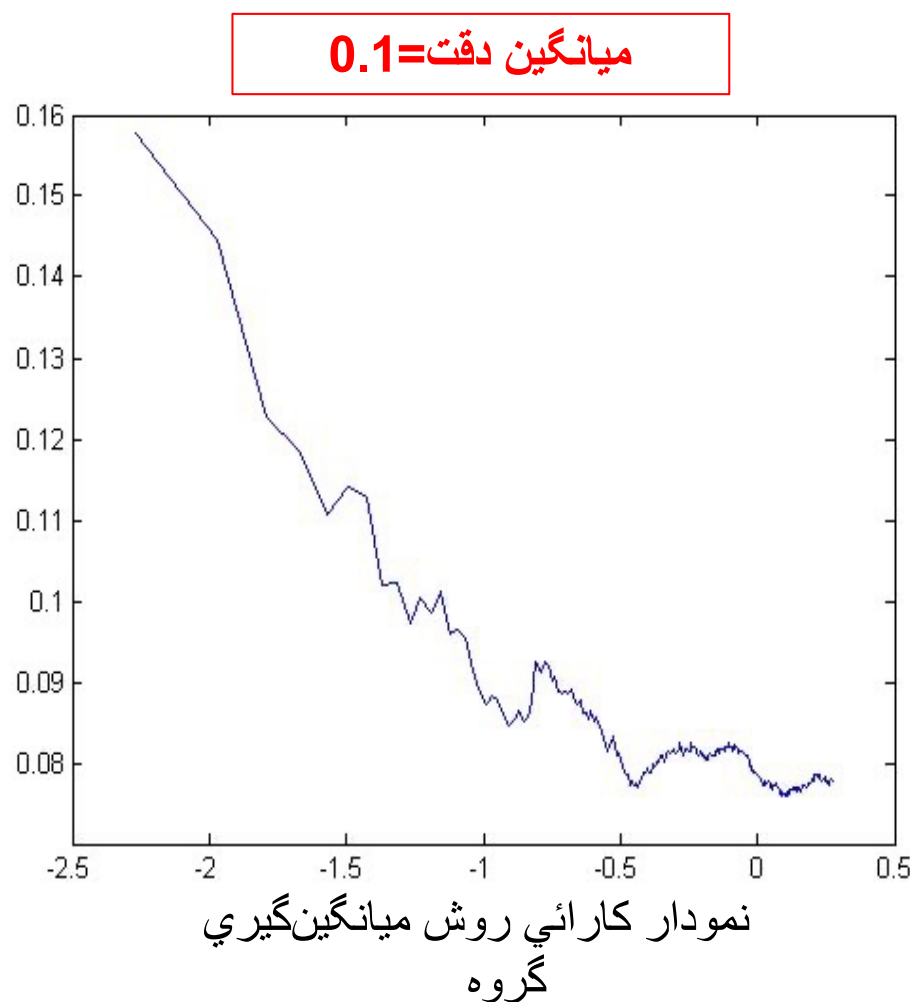
میانگین دقت=0.05



نمودار کارائی روش  
پایه

# روش پیشنهادي-بررسي نتايج (ادامه)

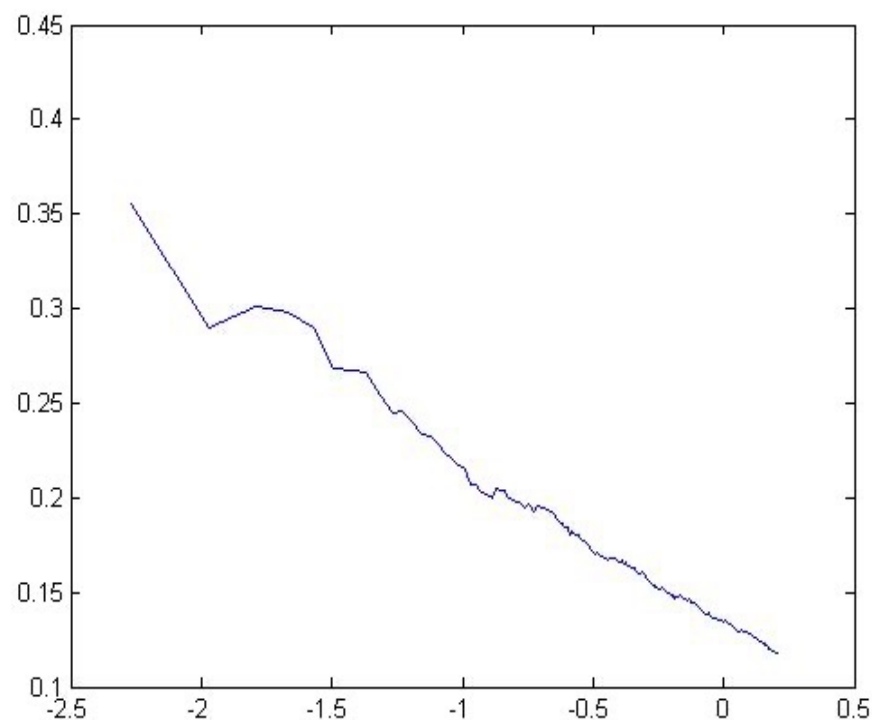
- نتیجه روش میانگین گیری گروه





# روش پیشنهادي-بررسي نتايج (ادامه)

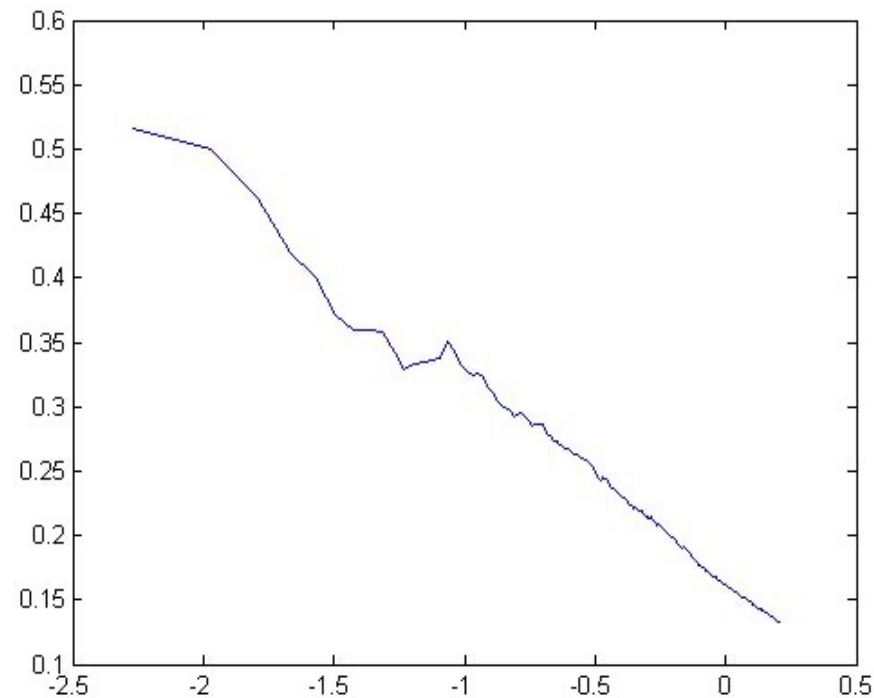
- نتايج روش پيشنهادي



نمودار کاري گروہ 1 بر روي زيرمجموعه تست داده‌هاي  
آموزشي

# روش پیشنهادي-بررسي نتايج (ادامه)

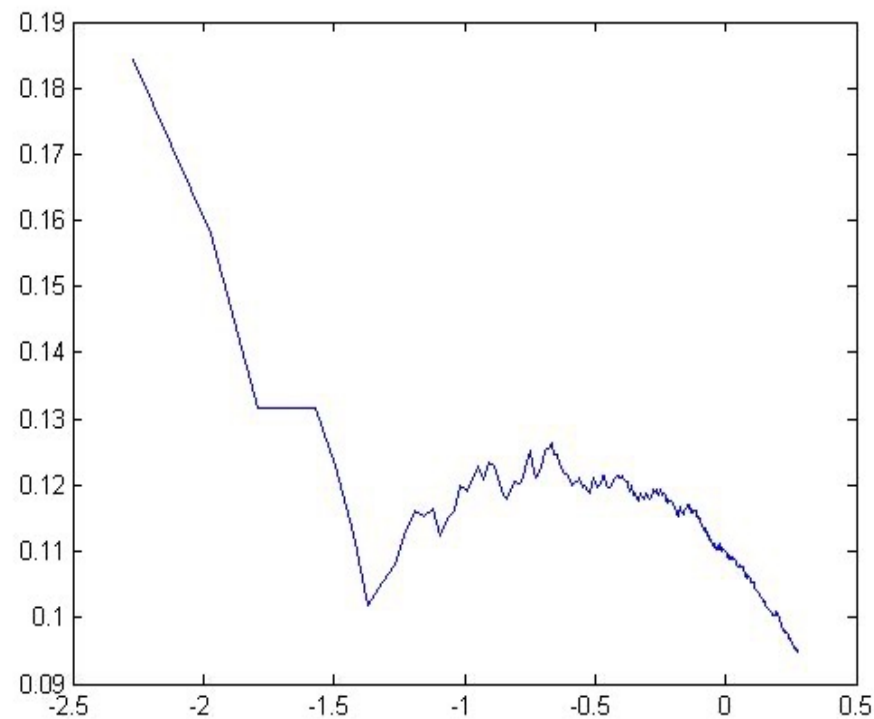
- نتايج روش پيشنهادي



نمودار کاري گروہ 2 بر روي زيرمجموعه تست داده‌هاي  
آموزشي

# روش پیشنهادي-بررسي نتايج (ادامه)

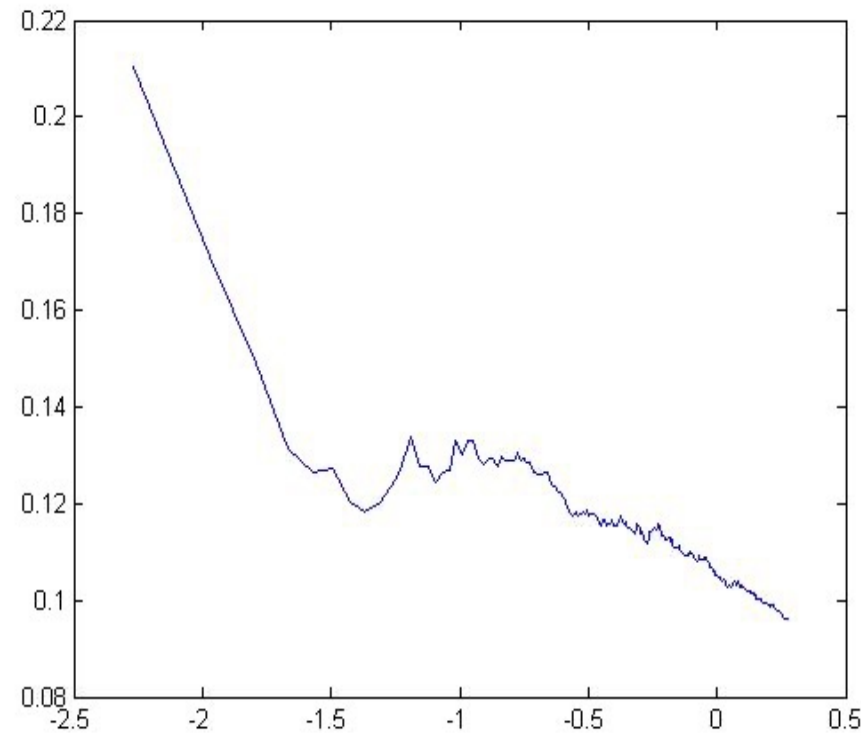
- نتايج روش پيشنهادي



نمودار کارائي گروه 2 با 5  
عضو

# روش پیشنهادي-بررسي نتايج (ادامه)

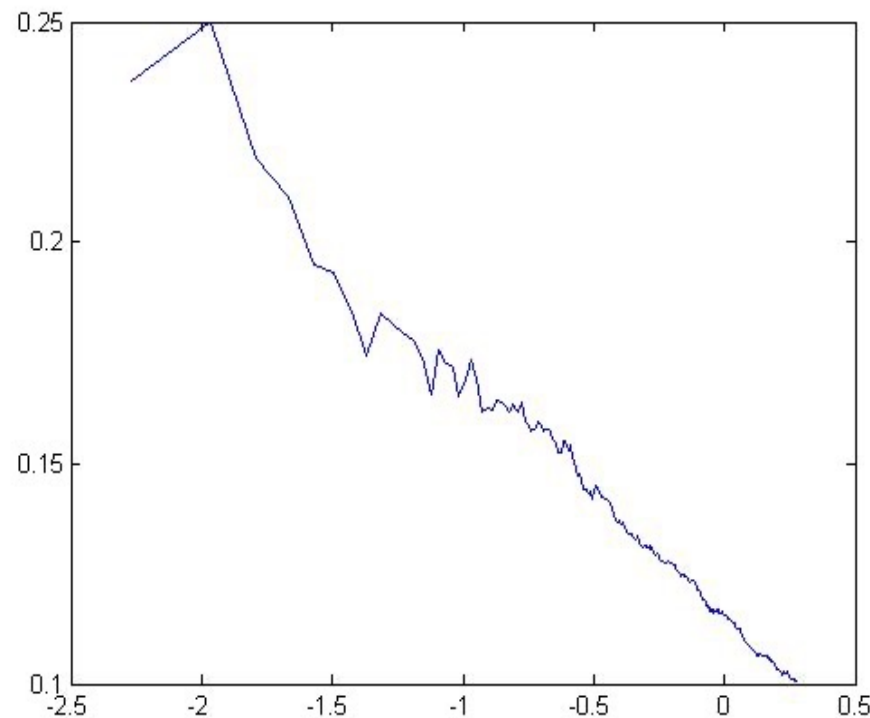
- نتايج روش پيشنهادي



نمودار کارائي گروه 2 با 20  
عضو

# روش پیشنهادي-بررسي نتايج (ادامه)

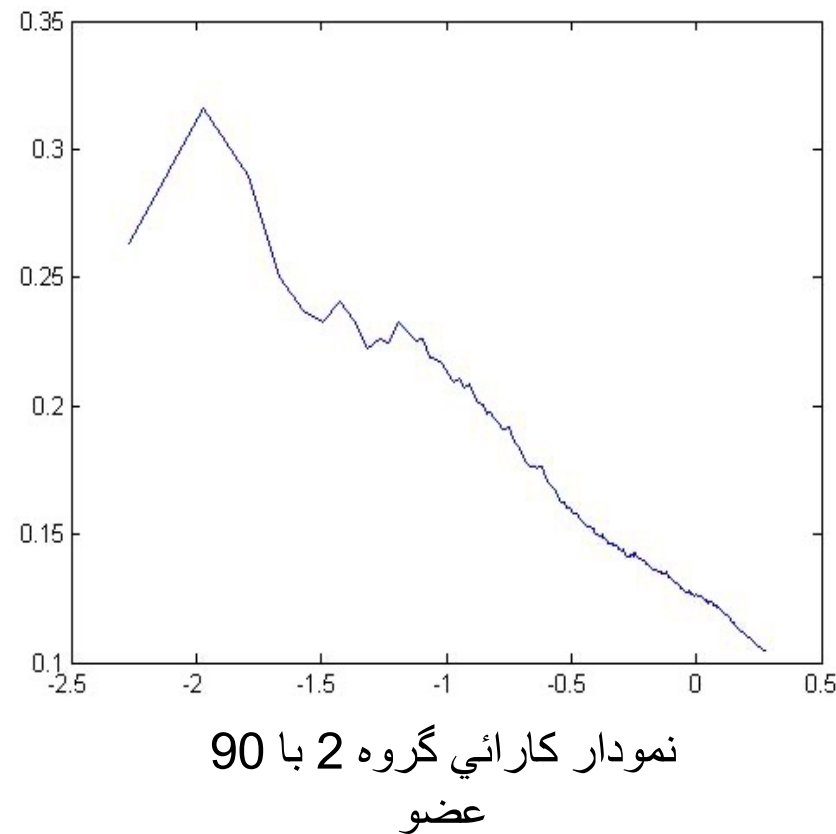
- نتايج روش پيشنهادي



نمودار کارائي گروه 2 با 50  
عضو

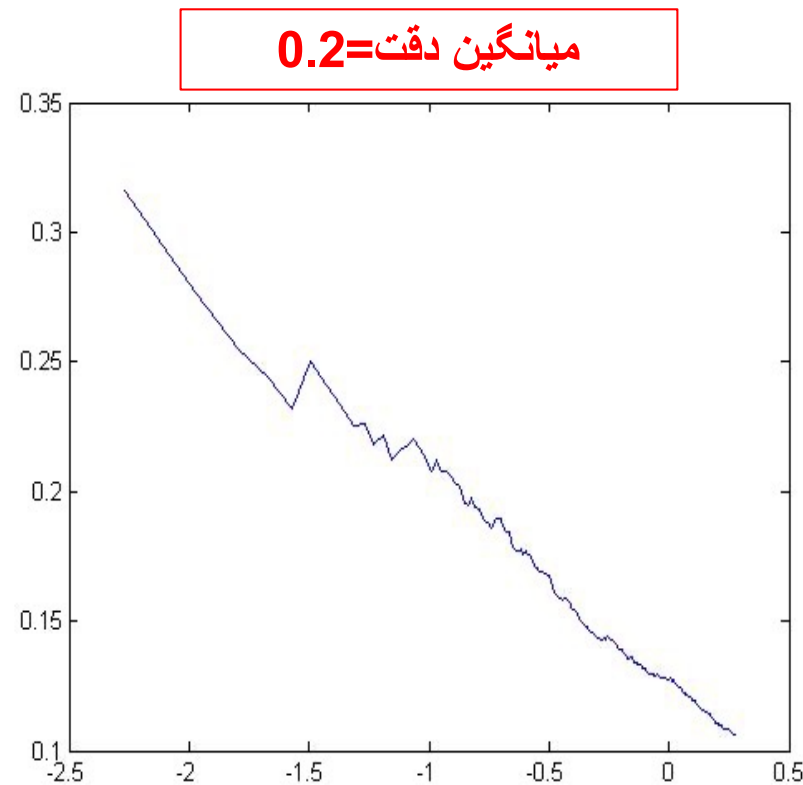
# روش پیشنهادي-بررسي نتايج (ادامه)

- نتايج روش پيشنهادي



# روش پیشنهادي-بررسي نتايج (ادامه)

- نتايج روش پيشنهادي



نمودار کارائي مدل

# روند ارائه مطالب

---

- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- پیش بینی نقشه تماس
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق



## نتیجه گیری

- ارائه روشی نوین پیش بینی نقشه تماس بر مبنای تکنیک ماشین گروهی
- نشان دادن کارایی روش پیشنهادی با بررسی نتایج
- مزایای روش پیشنهادی
  - کاهش خطا
  - کاهش واریانس
  - کاهش بیش پوشش
- معایب روش پیشنهادی
  - نیاز به حافظه بیشتر
  - افزایش هزینه آموزش
  - پیچیده تر شدن تحلیل سیستم

# روند ارائه مطالب

---

- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- پیش بینی نقشه تماس
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق

- پیش پردازش متفاوت برخی فایل های PDB، به جای حذف آن ها
- پیش پردازش فایل های MSA با دو آستانه شباهت متفاوت برای استخراج ویژگی های تک ستونی و جفت ستونی
- محاسبه ویژگی توزیع اسیدهای آمینه به روش تنظیم کننده ترکیبی دریکله، به جای محاسبه تقریبی احتمال رخداد آن ها
- استفاده از الفباها و روش های دیگر پیش بینی ساختار دوم
- استفاده از اندیس های AAindex3 به جای ویژگی تمایل

## پیشنهادهای (ادامه)

- استفاده از انواع دیگر ماشین گروهی به جای مدل پیشنهادی
- استفاده از مدل‌های یادگیری دیگر مانند ماشین بردار پشتیبان و مدل پنهان مارکوف، به جای شبکه عصبی
- استفاده همزمان از چند نوع مدل یادگیری در ماشین گروهی
- تعریف معیاری کارا تر برای تصمیم‌گیری درباره افزودن شبکه جدید به گروه
- حذف تماس‌های غیرمحمتمل از مجموعه تماس‌های پیش‌بینی شده، بر اساس قواعد بیوشیمی

# روند ارائه مطالب

---

- بیوانفورماتیک
- پروتئین
- مفاهیم محاسباتی پیش زمینه
- نقشه تماس پروتئین
- پیش بینی نقشه تماس
- روش پیشنهادی
- نتیجه گیری
- پیشنهادها
- دستاوردهای پژوهشی تحقیق

## Protein Contact Map Prediction Based On an Ensemble Learning Method

Narges Khatoun Habibi, Mohammad Hossein Saraei  
Advanced Database Systems, Data Mining and Bioinformatics Research Laboratory, Department of Electrical and  
Computer Engineering  
Isfahan University of Technology  
Isfahan, Iran, 84156-83111  
nhabibi@ec.iut.ac.ir, saraei@cc.iut.ac.ir

**Abstract**—Contact map is the simplified, 2D representation of protein spatial structure. Contact map prediction is an intermediate step to predict protein 3D structure. Ensemble learning-based model is a collection of learners that is more accurate than a single learner. In this paper we propose an ensemble learning method for contact map prediction. Results show that a considerable performance improvement is attainable using this approach instead of using a single model.

### I. INTRODUCTION

Proteins are the basic functional units of the cell which accomplish almost all the activities of life. The determination of the structure of proteins is an important step toward understanding the behavior of proteins. Experimental methods, such as x-ray crystallography and NMR techniques are not sufficient enough to allow for fast structural determination of the increasing number of newly discovered sequences. Therefore, employing computational techniques becomes vital.

Contact map is a simplified protein topology representation and it is an intermediate step to predict protein structure. It provides useful information about protein's structure. For example, secondary structure can easily be recognized from it. Alpha-helices appear as dense bands along the main diagonal since they involve contacts between one amino acid and its four successors, while Beta-sheets are dense bands parallel or anti-parallel to the main diagonal, etc [1].

Over the years, a variety of different approaches have been developed for contact map prediction, including statistical methods using correlated mutations [2,3], machine learning [1,4,5,6,7,8,9] and threading template-based voting [10].

The machine learning methods generate contact predictions by training the contact maps of known structures on a variety of sequence-based features including sequence profiles, secondary structure predictions, sequence conservation, correlated mutations and folding initiation

sites (i. e. I-sites). Different techniques of machine learning such as neural networks [4,5,10], support vector machines [6,7], evolutionary algorithms [1], and association rule based classification [9] have been used in this field. However, the accuracy of contact map prediction is still far from perfect.

In This paper we try to get use of a more advanced technique named ensemble learning. It has been shown that by using this technique, the overall model would have a better precision and performance in comparison to the single learner methods.

### II. DEFINITIONS

Background materials are presented in the following sub-sections.

#### A. Contact Map

Given a protein sequence  $P$  with  $N$  residues,  $P = (R_1, R_2, \dots, R_N)$ , two residues  $R_i$  and  $R_j$  ( $0 < i, j \leq N$ ), are considered to be in contact if their spatial distance is less than a threshold  $C_t$ . The distance between two residues, has various definitions in the literature. It could be defined as the distance between  $C_\alpha - C_\alpha$  atoms of the two residues [8],  $C_\beta - C_\beta$  [10], and the minimal distance between atoms belonging to the side chain or the backbone of the residue pair [11]. In this work, the distances between  $C_\beta - C_\beta$  are adopted. The distance function is [4]:

$$D_{ij} = |\vec{C}_{\beta,i} - \vec{C}_{\beta,j}| = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2} \quad (1)$$

Where  $\vec{C}_{\beta,i} = (x_i, y_i, z_i)$  is the three-dimensional coordination of the  $i^{th}$   $C_\beta$  atom (i.e. the  $C_\beta$  atom belonging to the  $i^{th}$  residue of the protein sequence).

We now could define a pairwise residue distance matrix  $D = \{D_{ij}\}$  for each protein, where each element  $C_{ij}$  of the contact map is set to 1 if the residue pair  $(R_i, R_j)$  is in contact, otherwise,  $C_{ij}$  is set to 0. In other words, the contact map of a protein with the length of  $N$  can be displayed by an  $N \times N$  matrix, whose elements can be defined as:

## Mining of Protein Primary Structure Data Using Committee Machines Approach to Predict Protein Contact Map

Narjes K. Habibi<sup>1</sup>, Kaveh Mahdavian<sup>1,2</sup> (IEEE Member), Mohammad H. Saraee<sup>1</sup>

<sup>1</sup>Advanced Database Systems, Data Mining and Bioinformatics Research Laboratory, Department of Electrical and Computer Engineering

Isfahan University of Technology (84156-83111)

<sup>2</sup>Isfahan Mathematics House

Isfahan, Iran

nhabibi@ec.iut.ac.ir, kaveh@ec.iut.ac.ir, saraee@cc.iut.ac.ir

**Abstract**—Committee machines approach has shown to be useful in different applications. Protein primary structure data contain valuable information to extract. In this paper we mine these data and predict protein contact map based on committee machines. Contact map is the simplified, two dimensional representation of protein spatial structure. Contact map prediction is of great interest due to its application in fold recognition and predicting protein tertiary structure. The results show that the performance of the committee is considerably better than a single model.

*artificial intelligence; machine learning; committee machine; ensemble learning; neural network; bioinformatics; protein contact map; contact map prediction*

### I. INTRODUCTION

Proteins are the basic functional units of the cell which accomplish almost all the activities of life. The determination of the structure of proteins is an important step toward understanding the behavior of proteins. Experimental methods, such as x-ray crystallography and NMR techniques are not sufficient enough to allow for fast structural determination of the increasing number of newly discovered sequences. Therefore, employing computational techniques becomes vital.

Contact map is a simplified protein topology representation, and it is an intermediate step to predict protein structure. It provides useful information about protein's structure. For example, secondary structure can easily be recognized from it. Alpha-helices appear as dense bands along the main diagonal since they involve contacts between one amino acid and its four successors, while Beta-sheets are dense bands parallel or anti-parallel to the main diagonal, etc [1].

Over the years, a variety of different approaches have been developed for contact map prediction, including statistical methods using correlated mutations [2,3], machine

learning [1,4,5,6,7,8,9] and threading template-based voting [10].

The machine learning methods generate contact predictions by training the contact maps of known structures on a variety of sequence-based features including sequence profiles, secondary structure predictions, sequence conservation, correlated mutations and folding initiation sites (i. e. I-sites). Different techniques of machine learning such as neural networks [4,5,10], support vector machines [6,7], evolutionary algorithms [1], and association rule based classification [9] have been used in this field. However, the accuracy of contact map prediction is still far from perfect.

In This paper we try to get use of a more advanced technique named committee machines. It has been shown that by using this technique, the overall model would have a better precision and performance in comparison to the single learner methods.

The rest of this paper is organized as follows. Section II provides some background information on contact maps, catastrophic forgetting and committee machines. In section III, the proposed model and applied input features are described in more details. Section IV presents experimental results. Finally, Section V provides some concluding remarks and future works.

### II. BACKGROUND MATERIALS

#### A. Contact Map

Given a protein sequence P with N residues,  $P = \{R_1, R_2, \dots, R_N\}$ , two residues  $R_i$  and  $R_j$  ( $0 < i, j \leq N$ ), are considered to be in contact if their spatial distance is less than a threshold  $C_T$ . The distance between two residues, has various definitions in the literature. It could be defined as the distance between  $C_\alpha - C_\alpha$  atoms of the two residues [8],  $C_\alpha - C_\beta$  [10], and the minimal distance between atoms belonging to the side chain or the backbone of the residue pair [11]. In this work, the distances between  $C_\alpha - C_\beta$  are adopted. The distance function is [4]:

# تشکر و قدردانی

- با سپاس از:
  - پروردگار یکتا
  - خانواده عزیزم
  - اساتید گرانقدر، آقای دکتر سرائی و آقای دکتر کربکندي
  - اساتید گرامی، آقای دکتر شیخ الاسلام و آقای دکتر موسوي
  - پروفیسور Kevin Karplus
  - آقای مهندس کاوه مهدویانی
  - تمامی معلمان و اساتید بزرگوارم در طول دوران تحصیل
  - اعضای آزمایشگاه داده کاوی و بیوانفورماتیک
  - دوستان خوبم
  - حضار محترم



