





دانشکده شیمی

رشته شیمی گرایش تجزیه

پایان نامه کارشناسی ارشد

مطالعه ارتباط کمی ساختار-فعالیت ترکیبات دارویی با استفاده از روش جنگل‌های

تصادفی

نگارنده: عطیه عسلی

اساتید راهنما:

دکتر ناصر گودرزی

دکتر منصور عرب چم جنگلی

شهریور ۱۳۹۵

اگر شایسته تقدیم باشد

تقدیم به:

«زیباترین نقش های عالم هستی»

پدر بزرگوارم؛

مادر عزیز و دلسوزم؛

و خواهر مهربانم

و

همه کسانی که دوستان دارم...

## پاسکزاری

سپاس بیکران ایزدمنان را که در طول این توفیق آمیز سخن میسر نمود تا منت پذیر آستان کبریایی اش باشم. امروز که به توفیق ایزد مهربان، راهی دیگر از زندگی را با موفقیت سپری کردم، پیشانی شکر بر سجده گاه عبودیت می سایم و بر خود واجب می دانم که از منت گذاران این راه قدر دانی نمایم و با شهادت قلم چند سطر بی ربه پاس زحمات بی دریغشان بنگارم. در همین راستا بر اساس روایت مشهور لم یسکر مخلوق لم یسکر الخالق، تحت سزاوار است نهایت سپاس قلبی خود را تقدیم حضور اساتید راهبانی کرامیم جناب آقای دکتر ناصر کوردزی و جناب آقای دکتر منصور عرب جم جفلی کردانم که زحمات بی شائبه ای متحمل گشته اند و در تمامی این مدت با بردباری مرار باهنایی فرموده اند و بی شک انجام مراحل مختلف این پایان نامه بدون حمایت و پشتیبانی ایشان امکان پذیر نبود. از داوران محترم جناب آقای دکتر قدوسی باقریان و سرکار خانم دکتر مصدرا لامور به خاطر نظرات ارزنده و اصلاحات بجا و دلسوزانه شان ممنون و سپاسگزارم.

در این تلاش کوچک با تمام عشق و اشتیاق از پدر و مادر عزیزم تختین آموزگاران زندگیم شکر کرده و بوسه بردستان مردانه پدر و چستان دعاگویی مادرم می زنم، آنان که امروز من آرزوی دیر زستان بود و از خداوند منان می خواهم عمری بینزاید تا گوشه ای از زحماتشان را جبران کنم، و دوست همیشگی ام خواهر عزیزم، آن که همواره حامی و مشوقم بوده و بیامودن روزهای سخت و آسان زندگی ام بدون دعای خیر و برکت و جودش غیر ممکن بود. سر و جودشان، همیشه سرسبز و

استور باد.

## تعهد نامه

اینجانب عطیه عسلی دانشجوی دوره کارشناسی ارشد رشته شیمی تجزیه دانشکده شیمی دانشگاه شاهرود نویسنده پایان نامه مطالعه ارتباط کمی ساختار-فعالیت ترکیبات دارویی با روش جنگل تصادفی تحت راهنمایی دکتر ناصر گودرزی و دکتر منصور عرب چم جنگلی متعهد می شوم :

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است .
- در استفاده از نتایج پژوهشهای محققان دیگر به مرجع مورد استفاده استناد شده است .
- مطالب مندرج در پایان نامه تاکنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است .
- کلیه حقوق معنوی این اثر متعلق به دانشگاه صنعتی شاهرود می باشد و مقالات مستخرج با نام « دانشگاه صنعتی شاهرود » و یا « Shahrood University of Technology » به چاپ خواهد رسید .
- حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تأثیرگذار بوده اند در مقالات مستخرج از پایان نامه رعایت می گردد.
- در کلیه مراحل انجام این پایان نامه ، در مواردی که از موجود زنده ( یا بافتهای آنها ) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است .
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری ، ضوابط و اصول اخلاق انسانی رعایت شده است .

امضای دانشجو

تاریخ

### مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج ، کتاب ، برنامه های رایانه ای ، نرم افزار ها و تجهیزات ساخته شده است ) متعلق به دانشگاه صنعتی شاهرود می باشد . این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود .
- استفاده از اطلاعات و نتایج موجود در پایان نامه بدون ذکر مرجع مجاز نمی باشد .

## چکیده

در بخش اول این مطالعه، مدل‌های QSAR برای پیش‌بینی فعالیت ضد HIV مشتقات پیریدینون ساخته شدند. برای انتخاب توصیف‌کننده‌ها از روش سهم گروه و رگرسیون خطی مرحله‌ای استفاده شد. سپس توصیف‌کننده‌های مهم به عنوان ورودی برای ایجاد مدل‌های QSAR با استفاده از روش جنگل‌های تصادفی (RF) و شبکه عصبی مصنوعی (ANN) استفاده شد. سری داده‌ها برای الگوریتم جنگل تصادفی به دو سری آموزش (۴۳ ترکیب) و سری آزمون (۱۰ ترکیب) تقسیم شدند. پارامترهای موثر بر روش جنگل‌های تصادفی شامل تعداد مشاهدات باقی‌مانده در هر گره ( $N_{tree}$ ) و تعداد توصیف‌کننده‌های انتخاب شده در هر مرحله افزای ( $M_{try}$ ) با استفاده از سری آموزش بهینه گردید. همچنین برای روش شبکه عصبی مصنوعی سری داده‌ها به سه دسته آموزش (۳۷ ترکیب)، ارزیابی (۸ ترکیب) و آزمون (۸ ترکیب) تقسیم شدند. بعد از بهینه‌سازی پارامترهای موثر، برای سری ارزیابی، بهترین مدل انتخاب گردید. ارزیابی مدل‌های بدست آمده توسط پیش‌بینی فعالیت سری آزمون انجام گرفت و مجذور ضریب همبستگی برای مدل‌های RF و ANN به ترتیب ۰/۹۲۹ و ۰/۹۳۵ به دست آمد.

در بخش دوم این تحقیق، روش‌های RF و ANN برای مدل‌سازی و پیش‌بینی صحیح فعالیت ضد HIV یکسری از مشتقات تiazول استفاده شد. داده‌ها به دو سری آموزش و آزمون به ترتیب شامل ۱۸ و ۶ ترکیب تقسیم شدند. برای انتخاب مهم‌ترین توصیف‌کننده‌ها روش رگرسیون مرحله‌ای مورد استفاده قرار گرفت. ارزیابی مدل‌های بدست آمده توسط سری آزمون انجام گرفت. مجذور ضریب همبستگی برای سری آزمون برای مدل‌های RF و ANN به ترتیب برابر با ۰/۹۸۲ و ۰/۸۹۶ به دست آمد. نتایج به دست آمده، توانایی پیش‌بینی مناسب مدل‌های پیشنهادی را در ارزیابی فعالیت ضد HIV مشتقات فوق نشان می‌دهد.

کلمات کلیدی: QSAR، ضد HIV، شبکه عصبی مصنوعی، جنگل‌های تصادفی، سهم گروه، رگرسیون خطی مرحله‌ای

نتایج حاصل از این پایان‌نامه در پوستری تحت عنوان

**“Application of random forest (RF) for the modeling of Anti-HIV  
activity of some dibenzyle pyridinon derivatives, as HIV-1 integrase  
inhibitor, using group contribution descriptors”**

در پنجمین سمینار دوسالانه کمومتریکس ایران دانشگاه تهران در آذرماه ۱۳۹۴ ارائه گردید.

## فهرست مطالب

۱- فصل اول .....	۱
۱-۱- ایدز.....	۲
۲-۱- ساختمان ویروس HIV.....	۲
۳-۱- انواع HIV.....	۴
۴-۱- میزان پایداری ویروس HIV.....	۴
۵-۱- سیر بیماری و علایم آن.....	۵
۶-۱- تشخیص بیماری.....	۶
۷-۱- درمان.....	۶
۸-۱- هدف‌های استفاده از داروهای ضد HIV.....	۷
۹-۱- دلایل استفاده از روش‌های تئوری برای پیش‌بینی فعالیت داروها.....	۸
۱۰-۱- مروری بر کارهای انجام شده.....	۸
۱۱-۱- ضرورت تحقیق.....	۱۱
۲- فصل دوم.....	۱۳
۱-۲- کمومتریکس.....	۱۴
۲-۲- ارتباط کمی ساختار- فعالیت (QSAR).....	۱۵
۳-۲- کاربردهای QSAR در طراحی دارو.....	۱۶
۴-۲- مراحل QSAR.....	۱۷
۵-۲- جمع آوری و پیش پردازش داده‌ها.....	۱۷
۶-۲- بهینه سازی ساختار مولکول‌ها.....	۱۸
۷-۲- محاسبه و غربالگری توصیف کننده‌ها.....	۲۰
۸-۲- توصیف کننده‌های سهم گروه.....	۲۲



- ۲-۹- بدست آوردن بهترین توصیف کننده‌ها برای ورود به مدل ..... ۲۳
- ۲-۹-۱- انتخاب توصیف کننده به روش رگرسیون مرحله‌ای ..... ۲۴
- ۲-۹-۲- انتخاب توصیف کننده بر مبنای روش سهم گروه ..... ۲۴
- ۲-۱۰- تقسیم‌بندی داده‌ها ..... ۲۵
- ۲-۱۱- مدل‌سازی ..... ۲۵
- ۲-۱۲- مقدمه ای بر شبکه عصبی مصنوعی ..... ۲۶
- ۲-۱۲-۱- ساختار و عملکرد نرون مصنوعی ..... ۲۷
- ۲-۱۲-۲- تابع انتقال ..... ۲۸
- ۲-۱۲-۳- ساختارهای شبکه ..... ۲۹
- ۲-۱۲-۴- آموزش شبکه‌های جلو سو با تکنیک پس انتشار ..... ۳۰
- ۲-۱۳- جنگل‌های تصادفی (RF) ..... ۳۱
- ۲-۱۳-۱- روش درخت رگرسیونی (تصمیم) ..... ۳۱
- ۲-۱۳-۲- الگوریتم تشکیل درخت رگرسیونی ..... ۳۵
- ۲-۱۳-۳- اندازه درخت و هرس کردن ..... ۳۷
- ۲-۱۳-۴- معرفی روش جنگل‌های تصادفی ..... ۳۸
- ۲-۱۴- ارزیابی قدرت پیش‌بینی مدل ..... ۴۳
- ۲-۱۴-۱- استفاده از پارامترهای آماری ..... ۴۳
- ۲-۱۴-۲- استفاده از نمودار برگشتی ..... ۴۶
- ۲-۱۴-۳- استفاده از نمودار خطای باقیمانده ..... ۴۷
- ۲-۱۴-۴- استفاده از سری تست ..... ۴۷
- ۲-۱۴-۵- استفاده از روش رد مرحله‌ای تکتک ..... ۴۷
- ۳- فصل سوم ..... ۴۹
- ۳-۱- مدل‌سازی فعالیت دارویی مشتقات پریدینون ..... ۵۰

- ۵۰-۲-۳- سری داده‌ها.....
- ۵۳-۳-۳- نرم‌افزارهای مورد استفاده.....
- ۵۳-۳-۳-۱- بسته نرم‌افزاری Hyperchem.....
- ۵۳-۳-۳-۲- بسته نرم‌افزاری دراگون.....
- ۵۴-۳-۳-۳- بسته نرم‌افزاری SPSS.....
- ۵۴-۳-۳-۴- نرم‌افزار MATLAB.....
- ۵۵-۴-۳- بهینه‌سازی ساختمان هندسی مولکول‌ها و محاسبه توصیف‌کننده‌ها.....
- ۵۵-۵-۳- انتخاب بهترین توصیف‌کننده‌ها.....
- ۵۵-۵-۳-۱- معرفی توصیف‌کننده‌ها به روش سهم گروه (GC).....
- ۵۵-۵-۳-۲- توصیف‌کننده‌های حاصل از دراگون و انتخاب بهترین توصیف‌کننده‌ها با روش رگرسیون مرحله‌ای (SR).....
- ۵۹-۶-۳- مدل‌سازی جنگل‌های تصادفی.....
- ۶۱-۶-۳-۱- مدل‌سازی جنگل‌های تصادفی با توصیف‌کننده‌های سهم گروه (GC-RF).....
- ۶۱-۶-۳-۲- مدل‌سازی جنگل‌های تصادفی با توصیف‌کننده‌های انتخاب شده به وسیله روش رگرسیون مرحله‌ای (SR).....
- ۶۴-۶-۳-۳- ارزیابی مدل جنگل‌های تصادفی.....
- ۷۰-۶-۳-۴- انتخاب توصیف‌کننده‌های مهم با روش جنگل‌های تصادفی.....
- ۷۱-۷-۳- مدل‌سازی شبکه عصبی مصنوعی.....
- ۷۱-۷-۳-۱- مدل‌سازی شبکه عصبی مصنوعی با توصیف‌کننده‌های سهم گروه (GC-ANN).....
- ۷۱-۷-۳-۲- مدل‌سازی شبکه عصبی مصنوعی با توصیف‌کننده‌های انتخاب شده به وسیله روش رگرسیون مرحله‌ای (SR-ANN).....
- ۷۷-۷-۳-۳- ارزیابی مدل‌های شبکه عصبی مصنوعی.....
- ۸۳-۷-۳-۳- ارزیابی مدل‌های شبکه عصبی مصنوعی.....

- ۸۸-۳-۷-۴- مقایسه مدل‌های ارائه شده با استفاده از پارامترهای آماری.....
- ۹۰-۳-۸- بررسی ارتباط بین توصیف‌کننده‌های منتخب و فعالیت مورد نظر.....
- ۹۴-۳-۹- بررسی میزان مشارکت توصیف‌کننده‌های منتخب در شبکه عصبی.....
- ۹۶-۳-۱۰- نتیجه‌گیری نهایی.....
- ۹۷-۴- فصل چهارم.....
- ۹۸-۴-۱- مدل‌سازی فعالیت دارویی مشتقات تیزول.....
- ۹۸-۴-۲- سری داده‌ها.....
- ۱۰۱-۴-۳- بهینه‌سازی ساختمان هندسی مولکول‌ها و محاسبه توصیف‌کننده‌ها.....
- ۱۰۱-۴-۴- انتخاب بهترین توصیف‌کننده‌ها با روش رگرسیون مرحله‌ای (SR).....
- ۱۰۳-۴-۵- مدل‌سازی و بهینه‌سازی پارامترهای موثر بر روش جنگل‌های تصادفی.....
- ۱۰۵-۴-۶- ارزیابی مدل جنگل‌های تصادفی.....
- ۱۰۵-۴-۶-۱- ارزیابی مدل RF با استفاده از داده‌های سری آزمون.....
- ۱۰۶-۴-۶-۲- ارزیابی مدل RF با استفاده از نمودار خطای باقیمانده.....
- ۱۰۶-۴-۶-۳- ارزیابی مدل RF توسط روش رد مرحله‌ای تک‌تک.....
- ۱۰۹-۴-۷- مدل‌سازی و بهینه‌سازی پارامترهای موثر بر شبکه عصبی مصنوعی.....
- ۱۱۴-۴-۸- ارزیابی مدل شبکه عصبی مصنوعی.....
- ۱۱۴-۴-۸-۱- ارزیابی مدل ANN با استفاده از داده‌های سری آزمون.....
- ۱۱۵-۴-۸-۲- ارزیابی مدل ANN با استفاده از نمودار خطای باقیمانده.....
- ۱۱۵-۴-۸-۳- ارزیابی مدل ANN توسط روش رد مرحله‌ای تک‌تک.....
- ۱۱۹-۴-۸-۴- مقایسه مدل‌ها با استفاده از پارامترهای آماری.....
- ۱۲۰-۴-۹- بررسی ارتباط بین توصیف‌کننده‌های منتخب و خاصیت مورد نظر.....
- ۱۲۴-۴-۱۰- بررسی میزان مشارکت توصیف‌کننده‌های منتخب در شبکه عصبی.....
- ۱۲۵-۴-۱۱- نتیجه‌گیری نهایی.....

پیشنهادات..... ۱۲۶

فهرست منابع..... ۱۲۷

## فهرست اشکال

- شکل ۱-۱: ساختمان ویروس HIV ..... ۴
- شکل ۱-۲: ساختمان یک نرون محاسباتی [۲۸] ..... ۲۷
- شکل ۲-۲: توابع انتقال [۳۰] ..... ۲۹
- شکل ۳-۲: شکل (الف) افراز قابل قبول و شکل (ب) افراز غیرقابل قبول ..... ۳۲
- شکل ۴-۲: نمودار درختی فضای افراز شده در افراز قابل قبول شکل (۲-۳) ..... ۳۳
- شکل ۵-۲: (الف) تقسیم بندی فضا، (ب) نمودار درختی تقسیم بندی فضا [۳۳] ..... ۳۳
- شکل ۶-۲: نمودار پراکنش داده‌های دو متغیره ..... ۳۵
- شکل ۷-۲: چهار افراز ممکن در راستای متغیر  $x_2$  ..... ۳۶
- شکل ۸-۲: سه افراز ممکن در راستای متغیر  $x_1$  ..... ۳۶
- شکل ۱-۳: اسکلت اصلی ترکیب مورد مطالعه ..... ۵۱
- شکل ۲-۳: نمایش نوع موقعیت‌ها (جایگاه‌های) استخلاف شده ..... ۵۶
- شکل ۳-۳: بهینه سازی پارامترهای  $M_{try}$  و  $n_{tree}$  در  $Node\ size$  بهینه ..... ۶۱
- شکل ۴-۳: بهینه سازی پارامترهای  $M_{try}$  و  $n_{tree}$  در  $Node\ size$  بهینه ..... ۶۳
- شکل ۵-۳: نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی توسط مدل‌های GC-RF و SR-RF برای داده‌های سری آزمون ..... ۶۵
- شکل ۶-۳: نمودار خطای باقیمانده برحسب مقادیر تجربی برای داده‌های سری آزمون با مدل‌های جنگل تصادفی GC-RF و SR-RF ..... ۶۵
- شکل ۷-۳: نمودار مقادیر پیش‌بینی شده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک توسط مدل‌های جنگل‌های تصادفی (GC-RF و SR-RF) بر حسب مقادیر تجربی ..... ۶۹
- شکل ۸-۳: نمودار مقادیر خطای باقیمانده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک توسط

- مدل‌های جنگل‌های تصادفی (SR-RF و GC-RF) بر حسب مقادیر تجربی ..... ۶۹
- شکل ۳-۹: نمودار اهمیت نسبی توصیف‌کننده‌های سهم گروه ..... ۷۰
- شکل ۳-۱۰: نمودارهای الف) تعداد گره‌ها، ب) تعداد دوره‌های آموزش و ج) مقادیر MSE برای تابع آموزش لونبرگ- ..... ۷۳
- شکل ۳-۱۱: نمودارهای الف) تعداد گره‌ها، ب) تعداد دوره‌های آموزش و ج) مقادیر MSE برای تابع آموزش لونبرگ- ..... ۷۴
- شکل ۳-۱۲: نمودارهای الف) تعداد گره‌ها، ب) تعداد دوره‌های آموزش و ج) مقادیر MSE برای تابع آموزش تنظیم بایزین به همراه تابع انتقال لگاریتم سیگموئید ..... ۷۵
- شکل ۳-۱۳: نمودارهای الف) تعداد گره‌ها، ب) تعداد دوره‌های آموزش و ج) مقادیر MSE برای تابع آموزش تنظیم بایزین به همراه تابع انتقال تانژانت سیگموئید ..... ۷۶
- شکل ۳-۱۴: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای تابع آموزش لونبرگ-مارکوات به همراه تابع انتقال لگاریتم سیگموئید ..... ۷۹
- شکل ۳-۱۵: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای تابع آموزش لونبرگ-مارکوات به همراه تابع انتقال تانژانت سیگموئید ..... ۸۰
- شکل ۳-۱۶: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای تابع آموزش تنظیم بایزین به همراه تابع انتقال لگاریتم سیگموئید ..... ۸۱
- شکل ۳-۱۷: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای تابع آموزش تنظیم بایزین به همراه تابع انتقال تانژانت سیگموئید ..... ۸۲
- شکل ۳-۱۸: نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی را برای داده‌های سری آزمون ..... ۸۴
- شکل ۳-۱۹: نمودار خطای باقیمانده بر حسب مقادیر تجربی برای داده‌های سری آزمون ..... ۸۴

- شکل ۳-۲۰: نمودار مقادیر پیش‌بینی شده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک توسط مدل شبکه عصبی مصنوعی بر حسب مقادیر تجربی ..... ۸۷
- شکل ۳-۲۱: نمودار مقادیر خطای باقیمانده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک توسط مدل شبکه عصبی مصنوعی بر حسب مقادیر تجربی ..... ۸۷
- شکل ۳-۲۲: مشارکت توصیف‌کننده‌ها در شبکه عصبی بهینه SR-ANN ..... ۹۵
- شکل ۴-۱: اسکلت اصلی شماره یک ..... ۹۸
- شکل ۴-۲: اسکلت اصلی شماره دو ..... ۱۰۰
- شکل ۴-۳: بهینه‌سازی پارمترهای  $M_{try}$  و  $n_{tree}$  در  $Node\ size$  بهینه ..... ۱۰۴
- شکل ۴-۴: نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی به وسیله داده‌های سری آزمون مدل جنگل‌های تصادفی ..... ۱۰۵
- شکل ۴-۵: نمودار خطای باقیمانده بر حسب مقادیر تجربی به وسیله داده‌های سری آزمون برای مدل جنگل‌های تصادفی ..... ۱۰۶
- شکل ۴-۶: نمودار مقادیر پیش‌بینی شده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک بر حسب مقادیر تجربی برای مدل جنگل‌های تصادفی ..... ۱۰۸
- شکل ۴-۷: نمودار مقادیر خطای باقیمانده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک بر حسب مقادیر تجربی برای مدل جنگل‌های تصادفی ..... ۱۰۸
- شکل ۴-۸: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای تابع آموزش لونیبرگ-مارکوات به همراه تابع انتقال لگاریتم سیگموئید ..... ۱۱۰
- شکل ۴-۹: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای تابع آموزش لونیبرگ-مارکوات به همراه تابع انتقال تانژانت سیگموئید ..... ۱۱۱
- شکل ۴-۱۰: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای تابع آموزش تنظیم‌بایزین به همراه تابع انتقال لگاریتم سیگموئید ..... ۱۱۲

- شکل ۴-۱۱: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و  
 د) مقادیر MSE برای تابع آموزش تنظیم بایزین به همراه تابع انتقال تانزانت سیگموئید..... ۱۱۳
- شکل ۴-۱۲: نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی به وسیله داده‌های  
 سری آزمون برای مدل شبکه عصبی مصنوعی ..... ۱۱۵
- شکل ۴-۱۳: نمودار خطای باقیمانده برحسب مقادیر تجربی برای داده‌های سری آزمون برای  
 مدل شبکه عصبی مصنوعی ..... ۱۱۵
- شکل ۴-۱۴: نمودار مقادیر خطای باقیمانده برای کل داده‌ها به وسیله روش رد مرحله‌ای تک-  
 تک برحسب مقادیر تجربی مدل شبکه عصبی مصنوعی ..... ۱۱۸
- شکل ۴-۱۵: نمودار مقادیر پیش‌بینی شده برای کل داده‌ها به وسیله روش رد مرحله‌ای تک-  
 تک برحسب مقادیر تجربی برای مدل شبکه عصبی مصنوعی ..... ۱۱۸
- شکل ۴-۱۶: مشارکت توصیف‌کننده‌ها در شبکه عصبی بهینه SR-ANN ..... ۱۲۴



## فهرست جداول

- جدول ۱-۳: استخلاف‌های ترکیب مورد مطالعه ..... ۵۱
- جدول ۲-۳: مقادیر نسبت داده شده به توصیف کننده‌ها ..... ۵۶
- جدول ۳-۳: توصیف کننده‌های سهم گروه ..... ۵۷
- جدول ۴-۳: توصیف کننده‌های انتخاب شده توسط رگرسیون مرحله‌ای ..... ۶۰
- جدول ۵-۳: ماتریس همبستگی کل توصیف کننده‌های انتخاب شده توسط رگرسیون مرحله‌ای ..... ۶۰
- جدول ۶-۳: کمترین مقادیر MSE همراه با  $M_{try}$  و  $n_{tree}$  متناظر با آن‌ها ..... ۶۲
- جدول ۷-۳: کمترین مقادیر MSE همراه با  $M_{try}$  و  $n_{tree}$  متناظر با آن‌ها ..... ۶۳
- جدول ۸-۳: نتایج حاصل از ارزیابی مدل‌های جنگل تصادفی با استفاده از داده‌های سری آزمون ..... ۶۴
- جدول ۹-۳: نتایج حاصل از رد مرحله‌ای تک‌تک برای مدل‌های به دست آمده توسط روش جنگل‌های تصادفی ..... ۶۷
- جدول ۱۰-۳: توصیف کننده‌های برتر با استفاده از روش جنگل‌های تصادفی ..... ۷۱
- جدول ۱۱-۳: توابع و پارامترهای شبکه‌های بهینه GC-ANN ..... ۷۷
- جدول ۱۲-۳: توابع و پارامترهای شبکه‌های بهینه SR-ANN ..... ۸۲
- جدول ۱۳-۳: نتایج حاصل از ارزیابی مدل‌های شبکه عصبی مصنوعی با استفاده از داده‌های سری آزمون ..... ۸۳
- جدول ۱۴-۳: نتایج حاصل از ارزیابی مدل‌های (GC-ANN و SR-ANN) به دست آمده توسط روش رد مرحله‌ای تک‌تک ..... ۸۵

جدول ۳-۱۵: پارامترهای آماری محاسبه شده برای ارزیابی و مقایسه مدل‌های SR-ANN ،	۸۹
GC-RF ، GC-ANN و SR-RF	۸۹
جدول ۳-۱۶: اثر متوسط توصیف‌کننده‌های به دست آمده توسط روش رگرسیون مرحله‌ای. ۹۰	
جدول ۳-۱۷: نمایش ارتباط فعالیت با مقدار توصیف‌کننده MATS8m برای بعضی از	
ترکیبات مورد مطالعه	۹۲
جدول ۴-۱: ساختار و مقادیر تجربی $pIC_{50}$ مشتقات تiazول	۹۹
جدول ۴-۲: ساختار و مقادیر تجربی $pIC_{50}$ مشتقات تiazول	۱۰۰
جدول ۴-۳: توصیف‌کننده‌های انتخاب شده توسط رگرسیون مرحله‌ای	۱۰۲
جدول ۴-۴: ماتریس همبستگی کل توصیف‌کننده‌های انتخاب شده توسط رگرسیون مرحله‌ای	
	۱۰۳
جدول ۴-۵: کمترین مقادیر MSE همراه با $M_{try}$ و $N_{tree}$ متناظر با آن‌ها	۱۰۴
جدول ۴-۶: نتایج حاصل از ارزیابی مدل جنگل تصادفی به دست آمده با استفاده از	
داده‌های سری آزمون	۱۰۵
جدول ۴-۷: نتایج حاصل از رد مرحله‌ای تک‌تک برای مدل جنگل‌های تصادفی	۱۰۷
جدول ۴-۸: توابع و پارامترهای شبکه‌های بهینه SR-ANN	۱۱۴
جدول ۴-۹: نتایج حاصل از ارزیابی مدل شبکه عصبی مصنوعی با استفاده از داده‌های سری	
آزمون	۱۱۴
جدول ۴-۱۰: نتایج حاصل از رد مرحله‌ای تک‌تک برای مدل شبکه عصبی	۱۱۷
جدول ۴-۱۱: پارامترهای آماری برای مدل‌های طراحی شده برتر توسط RF و ANN	۱۱۹
جدول ۴-۱۲: اثر متوسط توصیف‌کننده‌های به دست آمده توسط روش رگرسیون مرحله‌ای	
	۱۲۰

جدول ۴-۱۳: نمایش ارتباط فعالیت با مقدار GATS6p..... ۱۲۲



۱- فصل اول

مقدمه

## ۱-۱- ایدز

سندرم نقص ایمنی اکتسابی<sup>۱</sup> (AIDS) یک بیماری عفونی کشنده است که اولین بار در سال ۱۹۸۱ شناسایی شد. ویروس ایدز سلول‌های ایمنی موجود در خون (لنفوسیت‌ها) و سلول‌های ایمنی موجود در بافت‌ها مانند مغز استخوان، طحال، کبد و گره‌های لنفاوی را درگیر می‌سازد. این سلول‌ها در تولید پادتن برای مقابله با بیماری‌ها و سرطان‌ها نقش دارند. در مجموع باید گفت که ایدز یک نوع نقص ایمنی ثانویه است که در سیر عفونت با ویروس ایدز ایجاد می‌شود. ویروس HIV از خانواده رترو-ویروس‌ها<sup>۲</sup> و از زیر خانواده لنتی ویروس‌ها<sup>۳</sup> می‌باشد. رتروویروس‌ها، گروهی از ویروس‌ها هستند که ژنوم آن‌ها از RNA تشکیل شده است. بنابراین برای تکثیر خود به آنزیمی به نام آنزیم نسخه بردار معکوس<sup>۴</sup> وابسته هستند که ژنوم RNA آن‌ها را به DNA نسخه برداری کند تا بعد بتواند آن را به کمک آنزیم اینتگراز<sup>۵</sup> وارد ژنوم میزبان کند و به این ترتیب امکان تکثیر ویروس به وجود آید. کلمه رترو (معکوس) در نام این ویروس‌ها به همین خاطر است چرا که معمولاً نسخه برداری از DNA به RNA انجام می‌شود اما در این دسته از ویروس‌ها عکس این عمل صورت می‌گیرد. رتروویروس‌ها باعث کاهش توانایی سیستم ایمنی بدن میزبان می‌شوند. علایمی که در بیماری ایدز مشاهده می‌شود مربوط به بیماری‌هایی است که در اثر نقص دستگاه ایمنی بدن تولید می‌شوند.

## ۱-۲- ساختمان ویروس HIV

یک ذره کاملاً کروی با قطر ۱۲۰-۱۰۰ نانومتر است که در میکروسکوپ الکترونی این ویروس به شکل دوازده وجهی است که با غشای دولایه لیپیدی پوشیده شده است. در سطح غشای لیپیدی

---

<sup>۱</sup> - Acquired Immunodeficiency Syndrom

<sup>۲</sup> - Retrovirus

<sup>۳</sup> - Lentivirus

<sup>۴</sup> - Reverse transcriptase

<sup>۵</sup> - Integrase

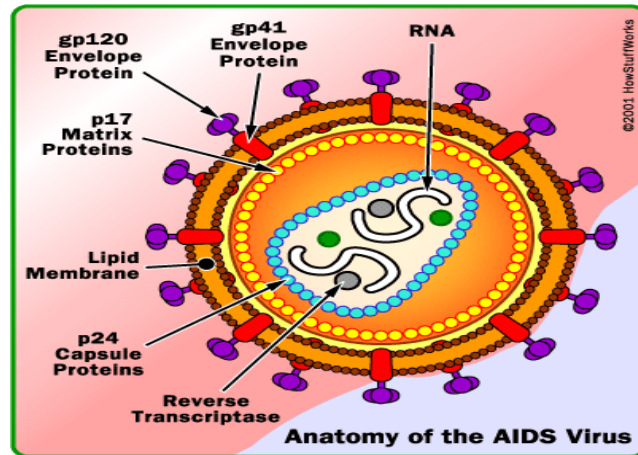
دو نوع گلیکوپروتئین (gp120، gp 41) وجود دارد که ویروس از طریق این گلیکوپروتئین‌ها قادر خواهد بود به گیرنده های  $CD_4^+$  موجود در سطح لنفوسیت T کمک کننده<sup>۱</sup> بچسبد. سطح داخلی غشای لیپیدی توسط ماتریکس احاطه شده و اطلاعات ژنتیکی مورد نیاز برای تکثیر ویروس را در خود جای داده است. این اطلاعات شامل دو کپی یکسان از ژنوم RNA ویروس و سه آنزیم اینتگرز، پروتئاز و نسخه بردار معکوس می‌باشد [۱]. ویریون<sup>۲</sup> (یک ذره ویروس که توان آلوده کردن سلول را دارد) از سطح سلول های آلوده جوانه زده و به محیط بیرون آزاد می‌شود و به همین علت است که ویروس HIV بطور مداوم پوشش خود را عوض می‌کند و به هنگام جوانه زدن از غشای سلول میزبان انواع متفاوتی پروتئین و گلیکوپروتئین جدا می‌کند که ترکیب آن منحصر به فرد می‌باشد. در شکل ۱-۱، به طور واضح لایه‌ها و هسته‌ی ویروس قابل مشاهده است. دو قطعه RAN به صورت جدا از هم در مرکز ویروس مستقر هستند. این ویروس‌ها حاوی آنزیم ترانس کریپتاز معکوس هستند و به کمک همین آنزیم است که پس از ورود به داخل سلول می‌توانند از روی RNA ژنوم دو رشته ای DNA را سنتز کنند. پس از این مرحله DNA سنتز شده، به هسته سلول میزبان رفته و توسط آنزیم‌های خود در DNA سلول میزبان رخنه می‌کند و سیستم همانند سازی کننده میزبان را وادار به همانند سازی از ژنوم خود می‌کند. در مراحل بعدی اجزای ژنوم ویروس HIV به سرعت سنتز شده و با پدیده جوانه زدن از سلول میزبان سبب آسیب‌های جدی به غشای سلول میزبان شده و در نهایت سلول تخریب می‌شود. با توجه به این که سلول‌های میزبان ویروس، گلبول‌های سفید هستند، این اتفاق سبب کاهش

---

<sup>۱</sup> - T helper lymphocyte

<sup>۲</sup> -Virion

شدید این سلول‌ها شده و موجب تضعیف و ایجاد نقص در سیستم ایمنی بدن می‌شود.



شکل ۱-۱: ساختمان ویروس HIV

### ۱-۳- انواع HIV

در حال حاضر دو گونه از HIV وجود دارد:

HIV-1 ✓

HIV-2 ✓

نوع غالب در جهان HIV-1 می‌باشد و هنگامی که بطور کلی و بدون مشخص کردن نوع درباره HIV بحث می‌شود، منظور HIV-1 می‌باشد. راه‌های انتقال هر دو نوع ویروس HIV یکسان است هرچند که ویروس HIV-2 مشکل‌تر از HIV-1 منتقل می‌شود و فاصله زمانی بین آلوده شدن با HIV-2 تا ایجاد بیماری طولانی‌تر است.

### ۱-۴- میزان پایداری ویروس HIV

ویروس HIV در اثر خشک شدن به سرعت غیرفعال می‌شود. در محیط کشت عاری از سلول و در درجه حرارت اتاق تا پانزده روز فعال بوده و در همان محیط در درجه حرارت ۳۷ درجه تا یازده روز زنده می‌ماند. در درجه حرارت ۶۰-۵۰ درجه سانتیگراد به مدت سی دقیقه عفونت زایی خود را از



دست می‌دهد. همچنین این ویروس در اثر استفاده از بعضی عوامل به مدت ده دقیقه در درجه حرارت اتاق به طور کامل غیر فعال می‌شود. این عوامل شامل هیپوکلریت سدیم ۰/۵٪، اتانول ۷۰٪، پارافرمالدهید ۰/۵٪، و قرار گرفتن در pH بیشتر از ۱۰ می‌باشد [۲]. از طرفی بدون درمان دارویی، بطور متوسط طی ۱۰ سال آلودگی به ویروس HIV به سمت ایدز پیش می‌رود که البته این مدت ۱۰ سال برای فردی است که تغذیه مناسبی دارد. اما فردی که در منطقه فقیر نشین است و بخوبی تغذیه نمی‌شود ممکن است بسیار سریعتر به سمت ایدز و نهایتاً مرگ پیش برود.

## ۱-۵- سیر بیماری و علایم آن

ابتلا به بیماری ایدز در اغلب موارد با عوارض ساده‌ای مانند اسهال، خارش پوست و افزایش خفیف حرارت بدن شروع می‌شود. این عوارض پس از چند هفته خودبخود برطرف شده و شخص خیال می‌کند که به سرماخوردگی دچار بوده است. از آنجا که ساخته شدن پادتن‌های موثر حدود ۱۴ روز طول می‌کشد در این زمان جدال بین پادتن‌ها و ویروس‌ها در جریان است. پس از سپری شدن این مرحله، دورانی فرا می‌رسد که حال شخص بیمار بهبود یافته اما در خون او پادتن ضد ایدز را می‌توان تشخیص داد و از خون وی ویروس ایدز را بدست آورد.

لنفوسیت T<sub>4</sub> که آلوده به ویروس شده است شروع به تقسیم شدن می‌کند و همین امر باعث تورم غدد لنفاوی می‌شود. تورم غدد لنفاوی گاه ماه‌ها یا سال‌ها بعد از ابتلای اولیه پدیدار شده و می‌تواند مدت‌ها به همان حال باقی بماند و در بسیاری از موارد ناراحتی‌های جدی دیگری به آن اضافه می‌شود. بیماران اکثراً به شدت لاغر می‌شوند چون دیواره روده آن‌ها مواد غذایی را بطور کامل جذب نمی‌کند. مرحله بعدی که در آن تصویر نهایی بیماری ایدز به نمایش گذاشته می‌شود حاصل از درهم شکسته شدن کامل سیستم دفاعی بدن است. تعداد سلول‌های لنفوسیت T<sub>4</sub> به شدت کاهش یافته و در این حالت انواع بیماری‌های باکتریایی، ویروسی و انگلی قادر هستند که فرد مبتلا را از پای درآورند.

## ۱-۶- تشخیص بیماری

تشخیص این بیماری از طریق آزمایش خون امکان پذیر است. سه نوع کلی آزمایش تشخیص

HIV وجود دارد که عبارتند از:

۱. تست آنتی بادی HIV<sup>۱</sup>: این تست نشان می دهد که فرد به HIV آلوده است یا نه.
۲. تست آنتی ژن P24<sup>۲</sup>: این تست بطور اولیه برای غربالگری نمونه های خون استفاده می شود ولی در برخی مناطق آن را بعنوان تست تشخیص HIV به کار می برند. آنتی ژن P24، یک پروتئین است که جزئی از ساختمان HIV بوده و در مراحل اولیه عفونت، به مقدار زیاد تولید می شود و بوسیله تست های تشخیصی می توان آن را در خون آشکار کرد. تست P24 می تواند آلودگی به HIV را قبل از تست آنتی بادی HIV آشکار کند. بنابراین، تست آنتی ژن P24 در تشخیص HIV در مراحل اولیه به کار می رود.
۳. تست میزان ویروس<sup>۳</sup>: این تست هنگامی استفاده می شود که شخص از آلوده بودن خود با HIV آگاه است و با این تست میزان ویروس در خون مشخص می شود.

## ۱-۷- درمان

در حال حاضر داروها در معالجه عفونت ایدز موثر نیستند. برای مقابله با عفونت ها یا پیشگیری

از آنها ممکن است آنتی بیوتیک تجویز شود. داروهای ضد HIV در سه دسته قرار می گیرند:

---

<sup>۱</sup> - HIV antibody test

<sup>۲</sup> - P24 Load test

<sup>۳</sup> - HIV Load test

۱) مهار کننده نوکلئوزیدی ترانس کریپتاز معکوس<sup>۱</sup>: زیدوودین<sup>۲</sup>، لامیوودین<sup>۳</sup>، دیدانوزین<sup>۴</sup>، زالسیتابین<sup>۵</sup> و استاوودین<sup>۶</sup> [۳].

۲) مهار کننده غیر نوکلئوزیدی ترانس کریپتاز معکوس<sup>۷</sup>: نونیراپینف دلاوردین و افایرنز [۴].

۳) مهار کننده پروتئاز: ایندیناویر، ساکویبی ناویر، نلفیناویر و آمپرنایویر.

این داروها ممکن است پیشرفت بیماری را به تاخیر اندازند و تجویز آنها باید توسط پزشکان با تجربه انجام پذیرد. با این داروها (مثل زیدوودین) می توان انتقال عفونت به نوزادان مادران باردار را کاهش داد. تحقیقات زیادی در رابطه با داروهای جدیدتر و واکسن آن از طریق مهندسی ژنتیک در دست انجام هستند.

## ۱-۸- هدف‌های استفاده از داروهای ضد HIV

۱. سرکوب ویروس و کاهش روند تکثیر آنها
۲. تحکیم و بهبود عملکرد سیستم ایمنی
۳. بهبود کیفیت زندگی بیماران با کاهش عوارض بیماری
۴. بهبود امید به زندگی و افزایش طول عمر بیماران
۵. کاهش احتمال سرایت

---

<sup>۱</sup> -Nucleoside Reverse Transcriptase Inhibitors (NRTIs)

<sup>۲</sup> -Zidovudine

<sup>۳</sup> -Lamivudine

<sup>۴</sup> -Didanosine

<sup>۵</sup> -Zalcitabine

<sup>۶</sup> -Stavudine

<sup>۷</sup> -Non-Nucleoside Reverse Transcriptase Inhibitors (NNRTIs)

## ۱-۹- دلایل استفاده از روش‌های تئوری برای پیش‌بینی فعالیت

### داروها

امروزه پیشرفت علم باعث بهبود روش‌های آزمایشگاهی گردیده است. اما به علت وجود محدودیت‌هایی همچون عدم قطعیت نتایج آزمایشگاهی، هزینه بر بودن، وقت گیر بودن و در دسترس نبودن امکانات، توسعه روش‌های نظری<sup>۱</sup> به عنوان یک ابزار کارآمد برای رفع این محدودیت‌ها مورد توجه قرار گرفته است. به‌کارگیری این روش‌ها می‌تواند علاوه بر پیش‌گویی فعالیت بیولوژیکی مورد نظر، به هدفمندتر شدن انجام آزمایش‌های تجربی و همین‌طور توضیح پارامترهای موثر بر نتایج این آزمایش‌ها و مکانیسم‌های درگیر کمک کند. از جمله روش‌های نظری، رابطه کمی ساختار - فعالیت<sup>۲</sup> (QSAR) است که در آن می‌توان از روش‌های مختلفی برای ایجاد روابط خطی و غیرخطی میان ساختار و خاصیت دارویی ترکیباتی که فعالیت دارویی آن‌ها با استفاده از روش‌های تجربی اندازه‌گیری شده‌اند، استفاده نمود.

## ۱-۱۰- مروری بر کارهای انجام شده

جلالی هروی و پرستار در سال ۲۰۰۰ در چارچوب مطالعات QSAR برای ایجاد مدل‌هایی که ارتباط دهنده‌ی فعالیت ضد HIV با ساختار ۱۰۷ مشتق (۱-۲) (هیدروکسی اتوکسی) متیل (۶-۶- فنیل تیو) تیمین (HEPT) باشد، مطالعاتی انجام دادند. در این کار مدل‌سازی توسط شبکه عصبی مصنوعی (ANN)<sup>۳</sup> نسبت به روش رگرسیون خطی چندگانه (MLR)<sup>۴</sup> نتایج بهتری را ارائه داد به

---

<sup>۱</sup> -Theoretical approaches

<sup>۲</sup> -Quantitative Structure- Activity Relationship

<sup>۳</sup> - Artificial Neural Network

<sup>۴</sup> - Multiple Linear Regression

طوری که میانگین مربع خطا (MSE)<sup>۱</sup> برای سری تست با استفاده از روش خطی چندگانه ۰/۷۸۰ و با استفاده از شبکه عصبی مصنوعی ۰/۳۷۲ به دست آمد [۵].

دوئلی<sup>۲</sup> و همکارانش در سال ۲۰۰۴ ارتباط کمی ساختار- فعالیت ۸۰ ترکیب از مشتقات HEPT که بازدارنده آنزیم نسخه بردار معکوس در ویروس HIV می‌باشند، را با استفاده از شبکه عصبی مصنوعی مورد بررسی قرار دادند و نتایج این مدل را با روش‌های رگرسیون خطی چندگانه (MLR) و حداقل مربعات جزئی (PLS)<sup>۳</sup> که توسط لوکو<sup>۴</sup> و فرتی<sup>۵</sup> ارائه شده بود، مقایسه کردند و از مقادیر مجذور ضریب همبستگی، قابلیت پیش‌بینی بهتر مدل ارائه شده توسط ANN را نتیجه گرفتند [۶].

عرب چم جنگلی در سال ۲۰۰۹ با استفاده از روش غیرخطی بر مبنای شبکه عصبی مصنوعی، مدلی مناسب برای فعالیت ضد HIV مشتقات پیرویل آریل سولفون‌ها ارائه داد به طوری که ضریب همبستگی برای سری تست ۰/۹۸۹ بود [۷].

در سال ۲۰۱۰ پاتل<sup>۶</sup> و همکارانش مطالعه QSAR را بر روی ۴۶ مشتق دی‌هیدرو-آلکوکسی-بنزیل-اکسوپیریمیدین (DABOs) که به عنوان بازدارنده‌های غیرنوکلئوزیدی ویروس HIV عمل می‌کنند، انجام دادند. مدل نهایی به دست آمده شامل ۴ توصیف‌کننده بود و با توجه به علامت ضرائب این توصیف‌کننده‌ها، نتیجه گرفتند که استخلاف‌های آگریز و استخلاف‌های با شاخه بزرگتر، تأثیر مثبتی بر بازدارندگی آنزیم نسخه‌بردار معکوس دارند [۸].

---

<sup>۱</sup> - Mean Square Error

<sup>۲</sup> - Douali

<sup>۳</sup> - Partial Least Square

<sup>۴</sup> - Luco

<sup>۵</sup> - Ferreti

<sup>۶</sup> - Patel

عرب چم جنگلی و اشرفی در سال ۲۰۱۳ با استفاده از روش غیرخطی بر مبنای شبکه عصبی مصنوعی، مدلی مناسب برای فعالیت مهارکنندگی نکرپتوز مشتقات [۳ و ۲ و ۱] تیادیازول و تیوفن ارائه دادند. قابلیت پیش‌بینی مدل ساخته شده که شامل ۷ توصیف‌کننده بود، توسط رد مرحله‌ای تک تک مورد ارزیابی قرار گرفت که ضریب همبستگی محاسبه شده برابر ۰/۹۴۴۶ بود [۹].

عرب چم جنگلی و عجم در سال ۲۰۱۵ با استفاده از توصیف‌کننده‌های سهم گروه و روش جنگل تصادفی، مدلی مناسب برای فعالیت بازدارنده‌های غیرنوکلئوزیدی مشتقات تیوکربامات ارائه دادند. قابلیت پیش‌بینی مدل ساخته شده که شامل ۱۲ توصیف‌کننده بود توسط رد مرحله‌ای تک تک و با استفاده از سری آزمون مورد ارزیابی قرار گرفت و مجذور ضریب همبستگی محاسبه شده برای پیش‌بینی خاصیت ضد HIV برابر ۰/۸۷۴۸ بود [۱۰].

## ۱۱-۱- ضرورت تحقیق

از آنجا که ایدز هم یکی از دلایل اصلی مرگ و میر در جهان است نیاز به درمان توسط مواد دارویی برای آن امری ضروری است. هرچند که جهش‌های<sup>۱</sup> مکرر ویروس HIV سبب افزایش مقاومت- دارویی<sup>۲</sup> و کاهش کارایی درمان توسط داروهای شیمیایی شده است، ولی نتایج تحقیقات نشان می‌دهد، اگر دو یا چند داروی ضد رتروویروسی همزمان مصرف شوند، مقاومت دارویی به میزان قابل توجهی کاهش می‌یابد. بنابراین درمان ترکیبی<sup>۳</sup> سبب گسترش تحقیقات جدید در جهت کشف داروهای گزینش پذیرتر و قویتر شده است و از آنجایی که بیشتر داروهای شیمیایی توسط آزمون و خطا سنتز می‌شوند بسیار وقت گیر و هزینه بر هستند [۱۱،۱۲].

در این پایان نامه سعی شده است با به کار گیری روش‌های QSAR مدل‌هایی از جمله روش جنگل‌های تصادفی<sup>۴</sup> و شبکه عصبی مصنوعی<sup>۵</sup> برای مدل‌سازی فعالیت دارویی و پیش‌بینی فعالیت بیولوژیکی داروهای ضد ایدز ارائه شود.

---

<sup>۱</sup> -Mutations

<sup>۲</sup> -Drug-Resistance

<sup>۳</sup> -Combination Therapy

<sup>۴</sup> -Rndom Forests

<sup>۵</sup> - Artificial Neural Networks





۲- فصل دوم

تئوری

## ۲-۱- کمومتریکس<sup>۱</sup>

کمومتریکس عبارت است از کاربرد روش‌های ریاضی و آمار به منظور برقراری ارتباط بین سنجش‌های انجام شده روی یک سیستم یا فرایند شیمیایی، طراحی عملیات آزمایشگاهی، دستیابی حداکثر به اطلاعات، انتخاب روش‌های بهینه اندازه‌گیری برای تفسیر و تحلیل بهتر اطلاعات شیمیایی [۱۳].

کمومتریکس شاخه‌ای نوپا در بین علوم می‌باشد. هرچند که نام آن از شیمی گرفته شده است ولی در بین شاخه‌های دیگر علوم هم دارای کاربردهای گسترده می‌باشد. کمومتریکس از آمار، جبرخطی و کامپیوتر برای آنالیز داده‌های شیمیایی استفاده می‌کند و هدف استخراج بیشترین اطلاعات از داده‌های شیمیایی می‌باشد. از نظر شیمی‌دانان فعالیت‌ها و خواص یک ترکیب ناشی از ویژگی‌های ساختاری آن است. هرگاه مطالعات به‌صورت ارتباط بین ساختارمولکولی و خواصی از نوع بیولوژیکی (مانند فعالیت دارویی) به کار گرفته شود، اشاره به ارتباط کمی ساختار-فعالیت (QSAR) خواهد داشت [۱۴].

در مطالعات QSAR سعی بر این است تا رابطه‌ی هماهنگ میان فعالیت‌های شیمیایی و فیزیکی با ویژگی‌های مولکولی پیدا شود، به گونه‌ای که بتوان این قواعد را برای ارزیابی فعالیت ترکیبات جدید به کار برد. در واقع نتایج این نوع مطالعات علاوه بر شفاف سازی نحوه‌ی ارتباط بین خواص مولکول‌ها و ویژگی‌های ساختاری آن‌ها، به پژوهشگران در پیش بینی رفتار مولکول‌های جدید براساس رفتار مولکول‌های مشابه کمک می‌کند [۱۵].

یکی از زمینه‌های مهم کاربرد کمومتریکس در مطالعاتی است که خواص مولکول‌ها را به

---

<sup>۱</sup> - Chemometrics

ویژگی‌های ساختاری آن‌ها نسبت می‌دهد. موارد خاصی از این تحقیقات و مطالعات شامل رابطه کمی ساختار-فعالیت (QSAR)<sup>۱</sup>، رابطه کمی ساختار-خصوصیت (QSPR)<sup>۲</sup>، رابطه کمی ساختار-سمیت<sup>۳</sup> (QSTR) و رابطه کمی ساختار-بازداری<sup>۴</sup> (QSRR) است. فرض اساسی در مطالعات QSAR این است که ساختار یک مولکول (یعنی خواص هندسی، الکترونی و چیدمان فضایی اتم‌ها) حاوی ویژگی‌هایی است که مسئول خواص فیزیکی، شیمیایی و زیستی ترکیبات شیمیایی می‌باشند. این ویژگی‌ها که به صورت کمی هم قابل نمایش هستند تحت عنوان توصیف کننده<sup>۵</sup> شناخته می‌شوند [۱۶].

برخی از کاربردهای کمومتریکس در شیمی تجزیه عبارتند از:

- (۱) شناسایی الگوها از داده‌های کروماتوگرافی
- (۲) پیش بینی خواص و فعالیت مولکول‌ها
- (۳) شناخت الگوها و طبقه بندی داده‌ها
- (۴) شبیه سازی و جستجوی طیف‌های جرمی
- (۵) پیش بینی میزان جا به جایی شیمیایی
- (۶) پیش بینی ضرایب گزینش پذیری الکترودهای یون‌گزین

## ۲-۲- ارتباط کمی ساختار - فعالیت (QSAR)

مدل‌سازی QSAR در حوزه سم‌شناسی متولد شد. در واقع تلاش‌ها برای کمی کردن روابط بین

---

<sup>۱</sup> -Quantitative Structure Activity Relationship

<sup>۲</sup> - Quantitative Structure Property Relationship

<sup>۳</sup> - Quantitative Structure Toxicity Relationship

<sup>۴</sup> - Quantitative Structure Retention Relationship

<sup>۵</sup> -Descriptor

ساختار شیمیایی و پتانسیل سمیت بحرانی قسمتی از دستاوردهای منتشر شده سم‌شناسی برای بیش از ۱۰۰ سال بود. در حال حاضر علم QSAR به شاخه‌های مختلف شیمی، بیوشیمی، شیمی پزشکی، داروسازی و شیمی زراعی راه یافته است و به ابزاری اساسی در طراحی دارو و آفت‌کش تبدیل شده و نقش رو به رشدی در علوم زیست محیطی دارد [۱۶].

مولکول‌های مشابه با تغییر کوچکی در ساختارشان می‌توانند فعالیت بیولوژیکی کاملاً متفاوتی داشته باشند [۱۷].

بررسی بیولوژیکی فرآورده‌های طبیعی و برخی از مشتقات آن‌ها که خواص دارویی دارند، انگیزه‌ای جهت پیدا کردن وابستگی‌های احتمالی ساختار شیمیایی با فعالیت بیولوژیکی ایجاد نمود و مشخص شد که برخی از واحدهای ساختار مولکول‌هایی که از نظر بیولوژیکی فعال هستند، در مولکول‌های دیگری که دارای همان ویژگی بیولوژیکی می‌باشند، وجود دارد. این امر رهنمودی برای طراحی ساختار ترکیباتی با فعالیت بیشتر و در عین حال عملکردی اختصاصی‌تر بود.

در حال حاضر اصلاح و تغییر ساختار یک ترکیب الگو، عمده‌ترین روش و روند دستیابی به داروهای جدید می‌باشد. تغییرات ساختاری منجر به پیدایش ویژگی‌های جدید می‌شود که در نهایت موجب تغییر پخش دارو در سلول بافت، تغییر در کیفیت اثر متقابل با جایگاه‌های فعال آنزیم و گیرنده‌ها و در نهایت تغییر سرعت واکنش‌ها در جایگاه‌های مورد نظر می‌گردد [۱۸].

## ۲-۳- کاربردهای QSAR در طراحی دارو

QSAR به راحتی می‌تواند برای کمک به بهینه‌سازی سنتز داروها مورد استفاده قرار گیرد. تحقیقات برای پیشنهاد داروهای جدید در روش‌های داروسازی سنتی تا حد زیادی برپایه‌ی آزمون و خطا و استخراج موادی با خاصیت دارویی از طبیعت می‌باشد. پس از شناسایی اصول مولکولی حیات، دانشمندان به اطلاعاتی ارزشمند در مورد نحوه‌ی اثر مولکول‌های دارویی بر بدن دست یافتند. در این

تلاش‌ها مشخص شد که اثر دارویی هر ترکیب به ویژگی‌های ساختاری آن ترکیب وابسته است. هر یک از این ویژگی‌های ساختاری باعث بروز یکی از آثار دارویی مربوط به آن می‌شود. در علم داروسازی مدرن با تکیه بر روش‌های QSAR، تولید دارو به یک روند تحقیقاتی کاملاً جهت‌دار تبدیل شده است. برای انجام یک مطالعه QSAR ابتدا باید مقادیر فعالیت دارویی ترکیبات از طریق تجربی بدست آید و همچنین توصیف‌کننده‌های مناسب از طریق تجربی یا با محاسبه توسط کامپیوتر تولید شوند. مرحله بعد دستیابی به ارتباط ریاضی بین ویژگی‌های ساختاری و مقادیر فعالیت است.

## ۲-۴- مراحل QSAR

به طور کلی ایجاد یک مدل QSAR شامل ۵ مرحله اصلی زیر است:

- ۱) جمع‌آوری و پیش‌پردازش داده‌ها
- ۲) بهینه‌سازی ساختار مولکول‌ها
- ۳) محاسبه و غربال‌گری توصیف‌کننده‌ها و انتخاب متغیر
- ۴) مدل‌سازی
- ۵) ارزیابی مدل

## ۲-۵- جمع‌آوری و پیش‌پردازش داده‌ها

اولین مرحله مدل‌سازی جمع‌آوری و انتخاب یک سری مولکولی است که مقادیر تجربی فعالیت مورد نظر آن‌ها باید با صحت قابل قبول اندازه‌گیری شده باشد. فاکتور محدود کننده در توسعه مدل‌های QSAR در دسترس بودن داده‌های تجربی با کیفیت بالا است. در واقع باید اعتبار و صحت هر مدل QSAR به صحت تعیین تجربی فعالیت مورد نظر بستگی دارد. داده‌های تجربی می‌توانند در

---

<sup>۱</sup> -Preprocessing

آزمایشگاه تولید یا از مقادیر منتشر شده در مقالات استخراج شوند. لازم به ذکر است که فعالیت ترکیبات باید در شرایط آزمایشگاهی یکسان تعیین شده باشد. سری مولکولی ممکن است شامل یک سری متجانس از ترکیبات شیمیایی باشد یا طبقه‌ای از مواد شیمیایی با یک تنوع ساختاری موجه را در برگیرد. این تنوع می‌تواند منجر به عمومیت یافتن بیشتر مدل شود. برای اجتناب از تغییرپذیری اضافی درون داده‌ها به علت استفاده از منابع اطلاعاتی مختلف، داده‌ها باید فقط از یک منبع یا منابع مشابه استخراج شوند یا به عبارتی دیگر، پیوستگی داده‌ها باید رعایت شود. بعد از اینکه سری داده جمع آوری شد لازم است که روی آن‌ها پیش پردازش انجام شود. این مرحله وابسته به طبیعت، مقدار و کیفیت داده‌ها است. پیش پردازش داده‌ها شامل حذف داده‌های نامربوط، نامطمئن و با نویز زیاد است که می‌تواند تاثیر منفی روی کل فرآیند مدل‌سازی داشته باشد. این روش‌ها بستگی به اهداف مدل‌سازی و دانش اولیه درباره خصوصیات داده‌ها دارد.

در مدل‌سازی QSAR سری داده‌ها عموماً به دو دسته سری آموزشی<sup>۱</sup> و سری آزمون<sup>۲</sup> تقسیم می‌شود. عملیات مدل‌سازی بر روی سری آموزشی انجام می‌شود. از سری آزمون نیز جهت بررسی قدرت پیش‌بینی و اعتبار مدل استفاده می‌شود. سری آزمون به نحوی انتخاب می‌شود که نماینده جمعیت مولکول‌های سری آموزشی باشد.

## ۲-۶- بهینه سازی ساختار مولکول‌ها

کمومتریکس از روش‌های مفیدی که قادر به استخراج اطلاعات موجود در داده‌ها هستند استفاده می‌کند [۱۹]. بدین منظور از شیمی محاسباتی<sup>۳</sup> استفاده می‌شود. شیمی محاسباتی

---

<sup>۱</sup> -Training set

<sup>۲</sup> -Test set

<sup>۳</sup> -Computational chemistry

ساختارهای مولکولی را به صورت پارامترهای عددی معرفی و رفتار آن‌ها را با معادلات کوانتومی و فیزیک کلاسیک شبیه سازی می نماید. این امر به دانشمندان امکان می دهد که بتوانند از این طریق به اطلاعات مولکول از جمله ساختار هندسی، انرژی و خواص الکترونیکی و اثرات حلال دست پیدا کنند. در سال‌های اخیر استفاده از روش‌های محاسباتی در بین شیمی‌دانان رواج یافته است. در این روش‌ها می توان به راحتی محاسبات را انجام داد، بدون اینکه از اصول اولیه و روش محاسبه آگاهی دقیق داشت. روش‌های محاسباتی در شیمی به چندین دسته تقسیم می شوند. روشی که در این تحقیق به کار گرفته شده، روش  $AM1^1$  است که جز روش‌های محاسباتی نیمه تجربی  $^2$  می باشد. در روش‌های نیمه تجربی، که در برنامه‌هایی نظیر HyperChem وارد شده‌اند، محاسبات براساس مکانیک کوانتومی صورت می گیرد. در این روش‌ها یک سری از مقادیر تجربی نیز وارد شده‌اند و فقط الکترون‌های لایه ظرفیت در محاسبه در نظر گرفته می شوند. در این تحقیق از MNDO  $^3$  یا روش اصلاح شده صرفنظر از همپوشانی‌های تفاضلی استفاده شده است.

روش  $AM1$  همان روش MNDO اصلاح شده است که روی مولکول‌های قطبی نتایج خوبی می دهد ولی پیوند هیدروژنی را نمی تواند در نظر بگیرد. در این روش هسته و لایه‌های داخلی به شکل یک هسته مرکزی  $^4$  در نظر گرفته می شود و محاسبات روی الکترون‌های ظرفیت انجام می گردد [ ۱۴ و ۲۰].  $AM1$  معمولاً گرمای تشکیل و انرژی اکتیواسیون را بسیار دقیق تر از MNDO پیش بینی می کند. از جمله ویژگی بارز روش نیمه تجربی این است که زمان محاسبات در این روش‌ها کوتاه تر است.

---

<sup>۱</sup> -Austin Methods

<sup>۲</sup> -Semi-empirical methods

<sup>۳</sup> -Modified neglected of differential overlap

<sup>۴</sup> -Core

## ۲-۷- محاسبه و غربالگری توصیف کننده ها

توصیف کننده‌ها مقادیر عددی هستند که ویژگی‌های مختلف ساختاری و الکترونی مولکول‌ها را به طور کمی نشان می‌دهند. یعنی هر توصیف کننده بیانگر خصوصیت ویژه‌ای از مولکول است که ممکن است بر خاصیت مورد نظر تاثیر گذار باشد. توصیف کننده‌های مولکولی به دو دسته کلی تقسیم می‌شوند: توصیف کننده‌های حاصل از اندازه‌گیری‌های تجربی (مانند قطبش پذیری و ممان دوقطبی) و توصیف کننده‌های مولکولی نظری (تئوری) که از ساختار مولکول مشتق شده و می‌توانند به دسته‌های بیشتری مطابق با انواع مختلفی از نمایش مولکولی<sup>۱</sup> یا بر اساس ابعاد توصیف کننده مولکولی تقسیم شوند.

محاسبه توصیف کننده‌ها یکی از مهمترین مراحل مدل‌سازی و تعیین و انتخاب توصیف کننده‌ها می‌باشد. توصیف کننده‌های مولکولی نتیجه نهایی یک استدلال و روش ریاضی هستند که اطلاعات شیمیایی را به رمز تبدیل می‌کنند و آن‌ها را به صورت یک نماد نشان می‌دهند طوری که ارائه دهنده ویژگی‌های یک مولکول به صورت یک عدد مفید باشند. در مدل‌سازی QSAR از توصیف کننده‌هایی باید استفاده شود که در ارتباط با ویژگی مورد نظر مولکول بوده و الگوریتم محاسباتی آن‌ها تعریف واضحی داشته باشد و همچنین تفسیر آن‌ها واضح باشد. همچنین باید تغییرات تدریجی در ساختار مولکول به تغییرات تدریجی در مقادیر توصیف کننده منجر شده و محدود به طبقه بسیار کوچکی از ساختارهای مولکولی نباشد. به طور کلی توصیف کننده‌ها براساس چگونگی بیان خصوصیات مولکول به چند دسته کلی تقسیم می‌شوند که به تعدادی از آن‌ها در زیر اشاره می‌شود [۲۱]:

**توصیف کننده‌های ساختاری<sup>۲</sup>:** این سری توصیف کننده‌ها عمومی ترین توصیف کننده‌های

---

<sup>۱</sup> -Molecular representation

<sup>۲</sup> -Constitutional descriptors



مورد استفاده می‌باشند که ساختار ترکیب را بدون هیچگونه اطلاعاتی درباره هندسه مولکولی آن در اختیار قرار می‌دهند.

**توصیف کننده‌های توپولوژیکی<sup>۱</sup>:** جزء ساده‌ترین نوع توصیف کننده‌ها می‌باشند و به ساختار فضایی مولکولی ارتباطی نداشته و تنها به نوع اتم‌ها، نوع پیوندها و نحوه ارتباط اتم‌ها به یکدیگر وابسته اند این پارامترها را می‌توان بدون بهینه کردن ساختار مولکول محاسبه کرد.

**توصیف کننده‌های هندسی<sup>۲</sup>:** نسبت به نوع توپولوژیکی از اهمیت بیشتری برخوردارند. یکی از خصوصیات بارز این توصیف کننده‌ها، ارتباط آن‌ها با ساختار سه بعدی مولکول‌ها می‌باشد. بنابراین قبل از محاسبه این توصیف کننده‌ها باید پایدارترین ساختار فضایی مولکول تعیین شود.

**توصیف کننده‌های الکترونی<sup>۳</sup>:** بیانگر خواص الکترونی مولکول‌ها می‌باشند و مقادیر آن‌ها اطلاعات مفیدی راجع به برهمکنش‌های الکترونی بین مولکولی و درون مولکولی ارائه می‌دهند.

**توصیف کننده‌های فیزیکو...شیمیایی<sup>۴</sup>:** بیانگر بعضی از خواص فیزیکی و شیمیایی مولکول‌ها می‌باشند و به ساختار مولکول وابستگی شدیدی دارند که از جمله می‌توان به نقطه ذوب و نقطه جوش اشاره نمود.

توصیف کننده‌های تئوری بدون نیاز به داده‌های تجربی به کمک نرم افزارهای مناسب برای طیف وسیعی از ترکیبات قابل محاسبه هستند. بنابراین توصیف کننده‌های تئوری در این مطالعات بسیار با اهمیت می‌باشند. برای محاسبه توصیف کننده‌ها، ابتدا مولکول‌ها با استفاده از نرم افزارهایی

---

<sup>۱</sup> -Topological descriptors

<sup>۲</sup> -Geometrical descriptors

<sup>۳</sup> -Electronic descriptors

<sup>۴</sup> -Physico-Chemical descriptors

مانند Chemoffice, PC Model, Hyperchem رسم شده، ساختار مولکولها به وسیله روش مکانیک کوانتومی<sup>۱</sup> و یا روشهای نیمه تجربی<sup>۲</sup> بهینه می‌شوند. سپس خروجی این نرم افزارها که حاوی اطلاعات سه بعدی ساختارهای مولکولی مورد نظر است به عنوان ورودی نرم افزارهای مخصوص محاسبه توصیف کننده مورد استفاده قرار می‌گیرد.

در نوع دیگری از تقسیم بندی، توصیف کننده‌ها به دسته‌های صفر بعدی، یک بعدی، دوبعدی و سه بعدی طبقه بندی می‌شوند. برای توضیحات بیشتر در مورد این نوع دسته بندی به مرجع [۲۲] مراجعه شود.

## ۲-۸- توصیف کننده‌های سهم گروه<sup>۳</sup>

امروزه در طراحی صنعتی، روش‌های ساده، کارآمد و در عین حال معتبر برای پیش‌بینی خواص مواد از روی ساختار مولکولی آنها ضروری است. یکی از این روش‌ها، روش سهم گروه است که امروزه به‌طور گسترده‌ای در پیش‌بینی خواص فیزیکی و ترمودینامیکی مواد به کار می‌رود. این روش که به رابطه جمع‌پذیری گروه‌ها<sup>۴</sup> نیز معروف است، ارتباطی بین خاصیت یک ماده با ساختار شیمیایی آن برقرار می‌کند. مبنای بوجود آمدن روش سهم گروه‌ها این است که اگر چه هزاران ترکیب شیمیایی در دنیای علم و تکنولوژی وجود دارد اما تعداد گروه‌های ساختاری و عاملی سازنده‌ی همه‌ی این ترکیبات، بسیار کمتر از تعداد این ترکیبات می‌باشد. برای ایجاد توصیف کننده‌های سهم گروه دو مرحله وجود دارد: در مرحله اول خواص مواد شناخته شده به ساختار شیمیایی آنها مرتبط می‌گردد تا گروه‌های اصلی در این مواد شناسایی شده و مقادیر جمع پذیر مولی آنها مشخص شود. در مرحله دوم سهم

---

<sup>۱</sup> -Quantum mechanics

<sup>۲</sup> -Semi-empirical

<sup>۳</sup> - Group Contribution

<sup>۴</sup> -Group additivity relation

یک گروه مشخص در یک مولکول با تمام مولکول‌ها یکسان در نظر گرفته می‌شود. با این ایده، روش سهم گروه فرض می‌کند که خاصیت فیزیکی یا ترمودینامیکی هر ماده را می‌توان از مجموع سهم همه‌ی گروه‌های موجود در مولکول آن ماده در خاصیت مورد نظر، به دست آورد. به این ترتیب یک روش عملی برای پیش‌بینی خواص تعداد زیادی ماده بر حسب تعداد کمتری پارامتر به دست می‌آید که این پارامترها، سهم گروه‌های ساختاری را در خاصیت مورد نظر مشخص می‌کنند. اکثر این روش‌ها تجربی هستند اما اطلاعات نظری راجع به سهم گروه‌ها در خواص مواد می‌تواند به عنوان راهنما در توسعه‌ی روابط سهم گروه‌ها به کار گرفته شود. البته لازم به توضیح است که تعریف گروه‌های تشکیل دهنده‌ی هر ماده، یک مقوله اختیاری است [۲۳]. هرچه تمایز بین گروه‌ها بیشتر باشد صحت روش سهم گروه افزایش می‌یابد تا جایی که سرانجام هر ترکیب فقط شامل گروه‌های مربوط به خودش باشد. توصیف‌کننده‌های سهم گروه می‌تواند شامل قطعات اصلی مولکول، محل استخلاف‌ها، نوع استخلاف‌ها و... باشد. در این پژوهش نیز از توصیف‌کننده‌های سهم گروه استفاده شده است که از روی ساختار ترکیب و بدون نیاز به نرم افزار محاسبه شده‌اند.

## ۲-۹- بدست آوردن بهترین توصیف‌کننده‌ها برای ورود به مدل

یکی از مهمترین مراحل QSAR، انتخاب توصیف‌کننده‌های مناسب می‌باشد. زیرا توصیف‌کننده‌های نامناسب کار برازش و مدل‌سازی را طولانی می‌کنند و تأثیری در بهبود نتایج نخواهند داشت. به این منظور باید در انتخاب توصیف‌کننده‌ها دقت لازم به عمل آید و توصیف‌کننده‌هایی که بیشترین و نزدیکترین ارتباط را به پارامتر مورد نظر دارند، انتخاب شوند. از جمله روش‌های انتخاب

متغیر می‌توان به رگرسیون مرحله‌ای<sup>۱</sup>، سهم گروه<sup>۲</sup>، الگوریتم ژنتیک<sup>۳</sup> و روش جایگزینی<sup>۴</sup> اشاره کرد.

## ۲-۹-۱- انتخاب توصیف کننده به روش رگرسیون مرحله‌ای

یکی از روش‌های انتخاب توصیف کننده‌های مناسب روش رگرسیون مرحله‌ای است که با محاسبه ضرایب همبستگی<sup>۵</sup> مشخص می‌شوند. ضریب همبستگی میزان نزدیکی اطلاعات بدست آمده به خط رگرسیون را نشان داده که مقدار آن بین ۰ و ۱ تغییر می‌کند. "صفر" یعنی هیچ ارتباطی بین فعالیت و توصیف کننده وجود ندارد و "یک" یعنی این ارتباط کامل است. بنابراین توصیف کننده‌هایی با ضرایب همبستگی بزرگ در برازش منحنی برای ساختن مدل استفاده می‌شوند. در این روش ابتدا توصیف کننده‌ای که بیشترین همبستگی با متغیر وابسته (خاصیت مورد نظر) را دارد، وارد مدل شده و با ورود هر توصیف کننده جدید توصیف کننده‌های موجود در مدل بررسی شده و اگر هر کدام از آنها سطح معناداری خود را از دست داده باشد قبل از ورود توصیف کننده جدید از مدل خارج می‌شود [۲۴].

## ۲-۹-۲- انتخاب توصیف کننده بر مبنای روش سهم گروه

انتخاب توصیف کننده‌های مناسب برای برقراری ارتباط بین ساختار مولکول و فعالیت، در ساخت مدل اهمیت بسیاری دارد. انتخاب توصیف کننده‌ها بر مبنای روش سهم گروه‌ها، یک روش بسیار ساده است که در آن ابتدا باید گروه‌های ساختاری مولکول‌ها را مشخص نمود. سپس گروه‌های ساختاری در هر ترکیب به عنوان توصیف کننده‌های آن ترکیب انتخاب می‌شود، در صورتی که هر یک

---

<sup>۱</sup> -Stepwise regression

<sup>۲</sup>- Group Contribution

<sup>۳</sup>-Genetic Algorithm

<sup>۴</sup>-Replacment Method

<sup>۵</sup>-Correlation coefficient

از این گروه‌های ساختاری در ترکیبی وجود نداشته باشد توصیف کننده مربوط به آن صفر در نظر گرفته می‌شود.

## ۲-۱۰- تقسیم‌بندی داده‌ها

قبل از ساختن مدل، سری داده‌ها معمولاً به سه زیر مجموعه تقسیم می‌شوند. اولین زیر مجموعه سری آموزش<sup>۱</sup> می‌باشد که اکثریت داده‌ها را در بر گرفته و از آن برای ساخت مدل غیرخطی استفاده می‌شود. دومین زیر مجموعه، سری ارزیابی<sup>۲</sup> است که در روند آموزش دخالت ندارد و مدل‌های به دست آمده از سری آموزش توسط آن ارزیابی شده و به این ترتیب توصیف‌کننده‌های بهینه برای ساخت مدل غیرخطی مشخص می‌شود. زیر مجموعه‌ی سوم سری تست<sup>۳</sup> است که در طول فرایند مدل‌سازی هیچ‌گونه دخالتی نداشته و از آن برای مقایسه مدل‌های مختلف استفاده می‌شود.

## ۲-۱۱- مدل‌سازی

مدل، در واقع یک رابطه ریاضی است که بیان‌کننده‌ی رابطه بین متغیر وابسته و مستقل می‌باشد و به کمک آن می‌توان با داشتن مقادیر متغیرهای مستقل، متغیر وابسته را ارزیابی کرد. پس از انتخاب مناسب‌ترین توصیف‌کننده‌ها، با استفاده از روش‌های آماری مختلف به جستجوی مدل مناسبی پرداخته می‌شود که بتواند ارتباط بین توصیف‌کننده‌های انتخابی و پارامترهای مورد مدل‌سازی را به درستی بیان کند. برای مدل‌سازی از روش‌های گوناگون خطی و غیر خطی می‌توان استفاده کرد. در زیر به دسته بندی کلی این روش‌ها اشاره می‌شود:

---

<sup>۱</sup>- Training set

<sup>۲</sup>- Validation set

<sup>۳</sup>- Test set

## روش‌های آماری چند متغیره (روش‌های خطی)

رگرسیون خطی چندگانه<sup>۱</sup> (MLR)

حداقل مربعات جزئی<sup>۲</sup> (PLS)

رگرسیون اجزای اصلی<sup>۳</sup> (PCR)

## روش‌های غیر خطی

شبکه عصبی مصنوعی<sup>۴</sup> (ANN)

جنگل‌های تصادفی<sup>۵</sup> (RF)

در این پایان نامه از روش‌های غیرخطی برای مدل‌سازی استفاده شده است.

## ۲-۱۲- مقدمه ای بر شبکه عصبی مصنوعی

شبکه عصبی مصنوعی یک الگوی پردازش اطلاعات است که از سیستم‌های عصبی بیولوژیکی مانند مغز الهام می‌گیرد. این الگو از تعداد زیادی عناصر پردازش کننده که به هم متصل هستند تشکیل شده است که به طور هماهنگ برای حل مسائل ویژه عمل می‌کنند و مانند مغز انسان با مثال، آموزش می‌بینند. در مغز انسان یک نرون سیگنال‌ها را از سایر نرون‌ها<sup>۶</sup> توسط بخشی از ساختار خود به نام دندریت جمع‌آوری می‌کند. دندریت پس از دریافت اطلاعات، آن‌ها را به شکل سیگنال به هسته

---

<sup>۱</sup>-Multiple Linear Regression

<sup>۲</sup>-Partial Least Square

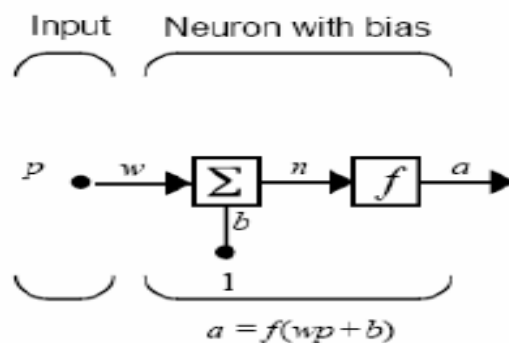
<sup>۳</sup>-Principle Component Regression

<sup>۴</sup>-Artificial Neural Networks

<sup>۵</sup> - Random Forests

<sup>۶</sup>-Neuron

سلول هدایت می‌کند. هسته پس از پردازش سیگنال‌های دریافتی، آن‌ها را از طریق جزء دیگری به نام اکسون به نرون‌های دیگر انتقال می‌دهد. این کار توسط سیناپس‌ها که ارتباط دهنده نرون‌ها هستند، صورت گرفته و بدین ترتیب فعالیت‌های مغزی انجام می‌شود [۲۶ و ۲۵]. با توجه به ساختار نرون محاسباتی در شکل ۱-۲ و مقایسه آن با نرون طبیعی، می‌توان ورودی‌ها را به دندریت، تابع محرک را به بدنه سلول، وزن‌ها را به شدت سیناپس‌ها و خروجی را به سیگنال گذرنده از اکسون تشبیه کرد [۲۷].



شکل ۱-۲: ساختمان یک نرون محاسباتی [۲۸]

## ۱-۱۲-۲- ساختار و عملکرد نرون مصنوعی

در شکل بالا کمیت‌های  $p$  و  $a$  به ترتیب ورودی و خروجی نرون می‌باشد. میزان تأثیر ورودی  $p$  روی خروجی  $a$  به وسیله پارامتر وزن ( $w$ ) تعیین می‌شود. ورودی دیگر یک مقدار ثابت  $1$  است که در جمله بایاس<sup>۱</sup> ضرب شده و سپس با  $wp$  جمع می‌شود. این حاصل جمع، ورودی خالص  $n$  برای تابع محرک  $f$  را تشکیل می‌دهد ( $a = f(n)$ ). در واقع  $b$  و  $w$  دو پارامتر تنظیم شونده در نرون‌ها می‌باشند و ایده اصلی شبکه عصبی این است که با تغییر مقادیر  $b$  و  $w$ ، شبکه یک رفتار یا تصمیم را اتخاذ کند

<sup>۱</sup>- Bias

## ۲-۱۲-۲- تابع انتقال<sup>۱</sup>

زمانی که ورودی‌ها با توجه به اهمیت آنها با یکدیگر جمع جبری شوند توسط تابع انتقال به نرون بعدی منتقل می‌شوند. در واقع تابع انتقال یکی از اجزای شبکه عصبی می‌باشد که برحسب کاربرد می‌تواند توابع متفاوتی باشد. دو تابع انتقال رایج در شبکه‌های عصبی که در این پژوهش برای بهینه‌سازی شبکه به کار گرفته شده است عبارتند از:

### - تابع انتقال لگاریتم سیگموئید<sup>۲</sup> ( $\log \text{sig}$ )

از این تابع انتقال در شبکه‌های پس انتشار استفاده می‌شود این تابع مقادیر ورودی را در محدوده منفی بی‌نهایت تا مثبت بی‌نهایت دریافت کرده و خروجی بین ۰ و ۱ تولید می‌نماید.

### - تابع انتقال تانژانت سیگموئید<sup>۳</sup> ( $\text{tansig}$ )

این تابع انتقال مقادیر ورودی را در محدوده منفی بی‌نهایت تا مثبت بی‌نهایت دریافت کرده و خروجی بین ۱+ و ۱- تولید می‌کند.

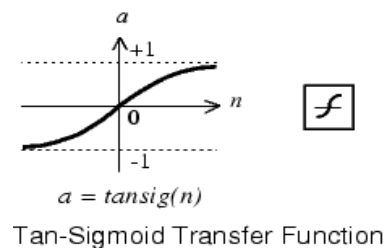
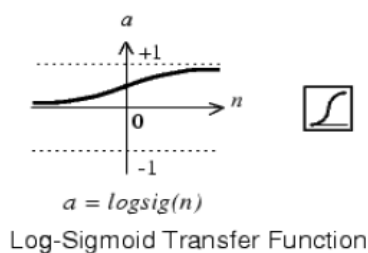
---

<sup>۱</sup> - Transfer function

<sup>۲</sup> - Logarithm sigmoid transfer function

<sup>۳</sup> - Hyperbolic tangent transfer function





شکل ۲-۲: توابع انتقال [۳۰]

## ۲-۱۲-۳- ساختارهای شبکه

بر اساس ساختار اتصال نرون‌ها، شبکه‌های عصبی مصنوعی را می‌توان در دو دسته‌ی کلی طبقه‌بندی کرد:

- شبکه‌های جلوسو<sup>۱</sup> که در آنها گراف حلقه ندارد.
  - شبکه‌های برگشتی<sup>۲</sup> که در گراف آنها بدلیل وجود بازخورد<sup>۳</sup> در ساختار شبکه، حلقه به وجود می‌آید.
- در معمول‌ترین خانواده‌ی شبکه‌های جلوسو که پرسپترون چندلایه<sup>۴</sup> نامیده می‌شوند، نرون‌ها در لایه‌هایی قرار می‌گیرند و اتصال بین آنها یک طرفه است. اتصالات مختلف سبب رفتارهای متفاوت شبکه‌ها می‌شود. بطور کلی می‌توان گفت شبکه‌های جلوسو، استاتیک هستند، به این معنی که از ورودی داده شده تنها یک دسته مقدار خروجی تولید می‌کنند نه یک دنباله از مقادیر خروجی! شبکه‌های جلوسو بی‌حافظه هستند و پاسخ آنها به یک ورودی مستقل از وضعیت قبلی شبکه است.

---

<sup>۱</sup>- Feed-Forward Networks

<sup>۲</sup>- Recurrent Networks

<sup>۳</sup>- Feedback

<sup>۴</sup>- Multi Layer Perceptron

از طرف دیگر شبکه‌های برگشتی، سیستم‌های دینامیک هستند و زمانی که یک دنباله‌ی ورودی جدید به آنها داده شود، خروجی نرون‌ها محاسبه می‌شود.

در اثر وجود مسی‌رهای بازخورد، ورودی هر نرون اصلاح می‌شود و این باعث می‌شود که شبکه وارد یک وضعیت جدید شود. ساختارهای مختلف شبکه، به الگوریتم یادگیری مناسب خود نیاز دارند.

## ۲-۱۲-۴- آموزش شبکه‌های جلو سو با تکنیک پس‌انتشار<sup>۱</sup>

پس‌انتشار خطا یک روش متداول آموزش با ناظر برای شبکه‌های جلوسو است یعنی برای به دست آوردن ارتباط بین متغیرهای ورودی و خروجی در یادگیری به الگوی آموزشی نیاز است. به طور کلی آموزش به کمک تکنیک پس‌انتشار بر طبق مراحل زیر انجام می‌شود [۳۱]:

۱- انتشار ورودی‌ها از نرون‌های ورودی به سمت نرون‌های خروجی

۲- اختصاص ماتریس وزن‌های تصادفی به هریک از اتصالات

۳- مقایسه خروجی‌های شبکه با مقادیر واقعی (مقادیر هدف) و محاسبه‌ی خطای شبکه

۴- پس‌انتشار خطا از نرون‌های خروجی به سمت نرون‌های ورودی و اصلاح وزن‌ها

۵- ارزیابی عملکرد شبکه با توجه به تابع کارآیی تعیین شده

مراحل فوق تا زمانی تکرار می‌شود که به حداکثر تکرار<sup>۲</sup> مجاز رسیده باشد یا مقدار تابع کارآیی از مقداری که تعیین شده کمتر باشد. شبکه عصبی انتخاب شده در این پژوهش یک شبکه جلوسو با

---

<sup>۱</sup> - Back propagation

<sup>۲</sup> - Epoch

الگوریتم آموزشی پس‌انتشار می‌باشد.

## ۲-۱۳- جنگل‌های تصادفی<sup>۱</sup> (RF)

یکی از روش‌های نوین مدل‌سازی و پیشرفته آماری، روش جنگل‌های تصادفی است که یکی از روش‌های یادگیری ماشین<sup>۲</sup> است. این روش در سال ۲۰۰۱ توسط بریمن<sup>۳</sup> متخصص آمار دانشگاه برکلی آمریکا و آدله کاتلر<sup>۴</sup> ارائه شد [۳۲]. قابلیت این روش، هنگامی که تعداد متغیرهای توضیحی<sup>۵</sup> زیاد است، بیشتر بروز می‌کند. بریمن روش درخت‌های رده بندی و رگرسیون (CART)<sup>۶</sup> را به عنوان جایگزینی برای روش‌های مبتنی بر شبکه معرفی کرده است. که به ترتیب برای طبقه بندی داده‌های کیفی و پیش بینی داده‌های کمی به کار می‌روند. در این پروژه هدف پیش بینی داده‌های کمی است لذا تنها به بخش رگرسیونی روش CART با نام درخت تصمیم اشاره می‌شود.

## ۲-۱۳-۱- روش درخت رگرسیونی<sup>۷</sup> (تصمیم)

این روش بر پایه تقسیم کردن مجموعه داده‌ها به قسمت‌های کوچکتر است. در روش رگرسیون خطی، پیش بینی کننده‌های رگرسیونی، مدل‌هایی هستند که در آن‌ها یک مدل پیش بینی واحد روی فضای کل داده‌ها انجام می‌دهد. اما ممکن است به دلیل تفاوت رفتار متغیر پاسخ در نواحی مختلف، برقراری یک مدل واحد، کارایی لازم را نداشته باشد. لذا یک روش جایگزین، تقسیم بندی

---

<sup>۱</sup> -Random Forests

<sup>۲</sup> -Machine learning

<sup>۳</sup> -Breiman

<sup>۴</sup> -Adele Cutler

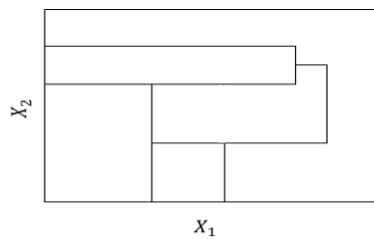
<sup>۵</sup> - Explanatory variable

<sup>۶</sup> -Classification And Regression Tree

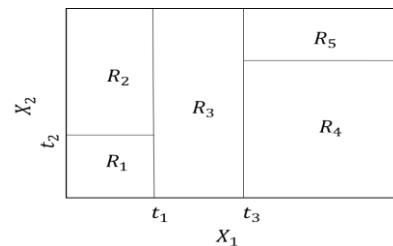
<sup>۷</sup> - Regression Tree

فضای داده‌ها به بخش‌های کوچکتر است تا بتوان رفتار متغیر پاسخ را به طور موضعی مدل سازی نمود. در روش درخت رگرسیونی هدف این است که مقادیر متغیرهای توضیحی ( ابر مکعب در فضای  $m$  بعدی) در هر ناحیه، به گونه‌ای تقسیم بندی شود که داده‌های واقع در هر بخش تا حد ممکن همگون باشند به طوری که بتوان آن‌ها را توسط ساده ترین مدل پیش بینی نمود.

به طور کلی، افراز<sup>۱</sup> فضای متغیرهای توضیحی را می‌توان به هر شکلی انجام داد. اما در روش درخت رگرسیونی، تنها افرازی قابل قبول است که تمامی نواحی ساخته شده به شکل مربع یا مستطیل باشند. در شکل ۲-۳ نمونه‌ای از افراز قابل قبول و غیرقابل قبول در فضای دومتغیره با روش درخت رگرسیونی نشان داده شده است.



(ب)

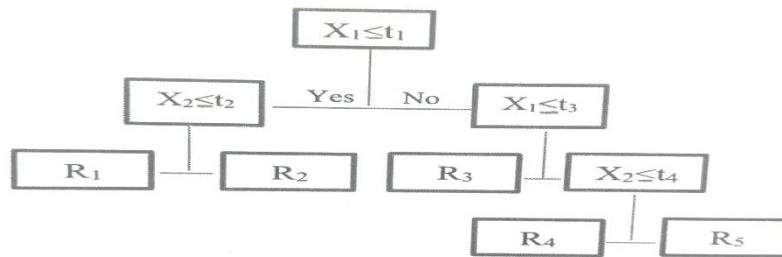


(الف)

شکل ۲-۳: شکل (الف) افراز قابل قبول و شکل (ب) افراز غیرقابل قبول

<sup>۱</sup>- Partition

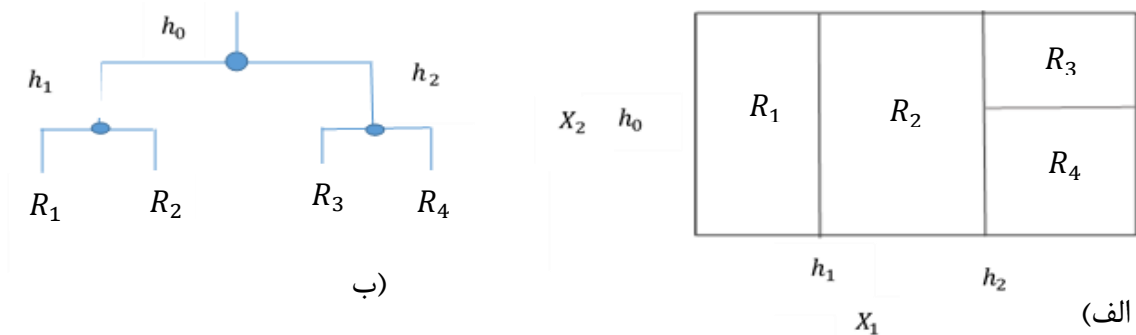
در شکل ۲-۴، نوع دیگری از نمایش این افراز نشان داده شده است. دلیل نامگذاری درخت رگرسیونی را می توان به نمایش درختی افراز فضای متغیرهای ورودی نسبت داد که در آن هر مشاهده از نقطه بالای نمودار درختی وارد شده و در یکی از نواحی افراز قرار می گیرد.



شکل ۲-۴: نمودار درختی فضای افراز شده در افراز قابل قبول شکل (۲-۳)

تقسیم بندی فضای مذکور توسط یک الگوریتم سلسله مراتبی انجام می شود و در آن به سوالات زیر پاسخ داده می شود:

- ۱) تقسیم فضا در راستای محور کدام یک از متغیرهای توضیحی و در چه مقداری از آن متغیر بایستی ایجاد شود؟
- ۲) معیار انتخاب راستای یاد شده چیست؟
- ۳) تقسیم فضا تا چه مرحله ای بایستی ادامه یابد؟



شکل ۲-۵: الف) تقسیم بندی فضا، ب) نمودار درختی تقسیم بندی فضا [۳۳]

به مقداری از هر متغیر که افراز در آن نقطه انجام می‌گیرد، گره<sup>۱</sup> می‌گویند. مثلاً در حالت دو متغیر توضیحی اگر تقسیم بندی فضا به صورت قسمت الف شکل ۲-۵ باشد آنگاه  $h_1, h_2, h$  مجموعه گره‌ها را تشکیل می‌دهد. اگر به ساختار درختی این تقسیم بندی (قسمت ب شکل ۲-۳) توجه کنیم، در واقع گره محلی است که زیر شاخه‌های درخت بهم متصل می‌شوند.

لازم به ذکر است که در روش درخت تصمیم، درخت‌ها دوتایی<sup>۲</sup> هستند به این معنی که هر زیر فضا (شاخه) فقط می‌تواند به دو مجموعه (شاخه) تقسیم شود.

اکنون با فرض اینکه فضای متغیرهای  $X_1$  و  $X_2$  به چهار ناحیه  $R_1, R_2, R_3, R_4$  تفکیک شده باشد، به طور قطع هر ناحیه  $m=1, \dots, 4$  شامل زیر مجموعه‌ای از داده‌های مشاهده شده به نام  $A_m$  است، به طوری که  $U_m^S = 1^{A_m} = S$ . که در آن فضای کل متغیرهای توضیحی می‌باشد. بنا به هدف دنبال شده در روش رگرسیونی، می‌توان رویه پاسخ را در ناحیه  $R_m$  توسط یک مدل ساده یعنی  $\hat{Y}_{R_m} = C_m$  تقریب نمود. بنابراین روش درخت رگرسیونی به صورت زیر به صورت زیر بیان می‌شود.

$$\hat{Y} = \sum_{m=1}^S c_m I_{A_m}(X_1, X_2) \quad (1-2)$$

$$I_{A_m}(X) = \begin{cases} 1, & X \in A_m \\ 0, & X \notin A_m \end{cases} \text{ که در آن } c_m \text{ها مقادیری ثابت می‌باشند و}$$

با توجه به روش کمترین توان‌های دوم خطا، به آسانی ثابت می‌شود که  $c_m$  میانگین مقادیر

متغیر پاسخ واقع در ناحیه  $R_m$  است. به عبارت دیگر

---

<sup>۱</sup> -Node

<sup>۲</sup> -Binary tree

$$\hat{y}_{R_m} = c_m = \bar{y}_{A_m} = \frac{1}{n_m} \sum_{y_j \in A_m} y_j \quad (2-2)$$

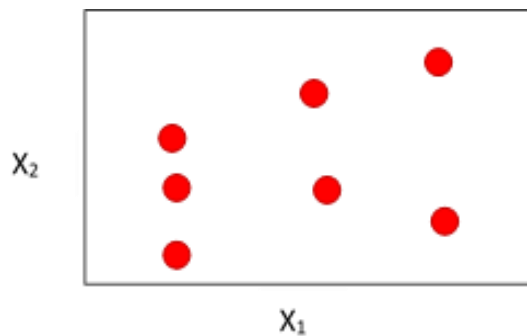
که در آن  $n_m$  تعداد مشاهدات مجموعه  $A_m$  است.

## ۲-۱۳-۲- الگوریتم تشکیل درخت رگرسیونی

تشکیل درخت رگرسیونی یا افراز فضای نمونه، طی سلسله مراتبی انجام می‌شود، که هر مرحله، براساس آن است که تفکیک فضا در راستای کدام متغیر و در چه مقدار از آن متغیر باید انجام گیرد. بدین منظور جزئیات الگوریتم را با یک مثال شرح می‌دهیم.

فرض کنید تعداد ۷ مشاهده برای دو متغیر  $X_1$  و  $X_2$  در اختیار باشد که نمودار پراکنش<sup>۱</sup> آن‌ها

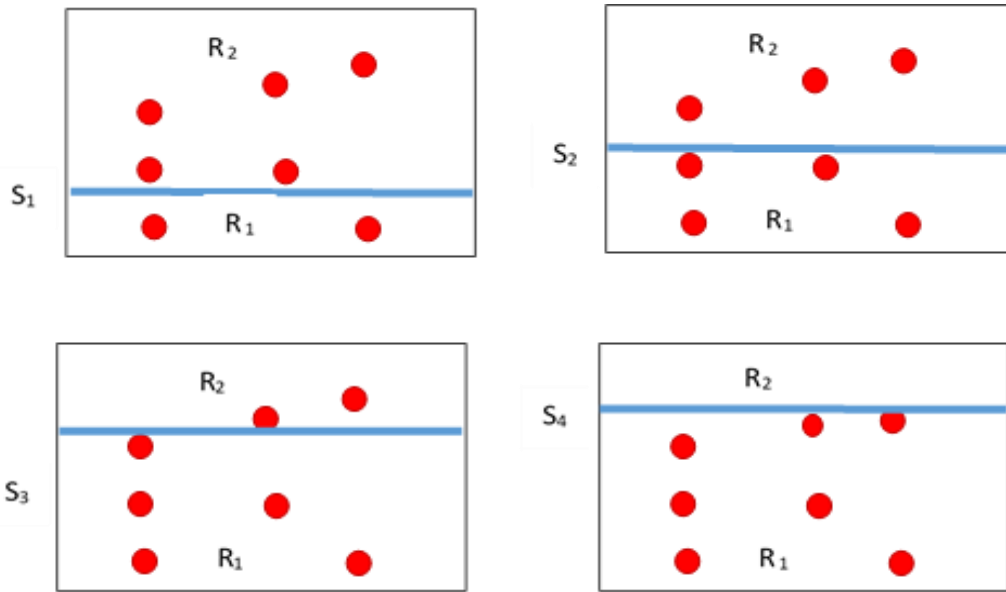
در شکل ۶-۲ آمده است.



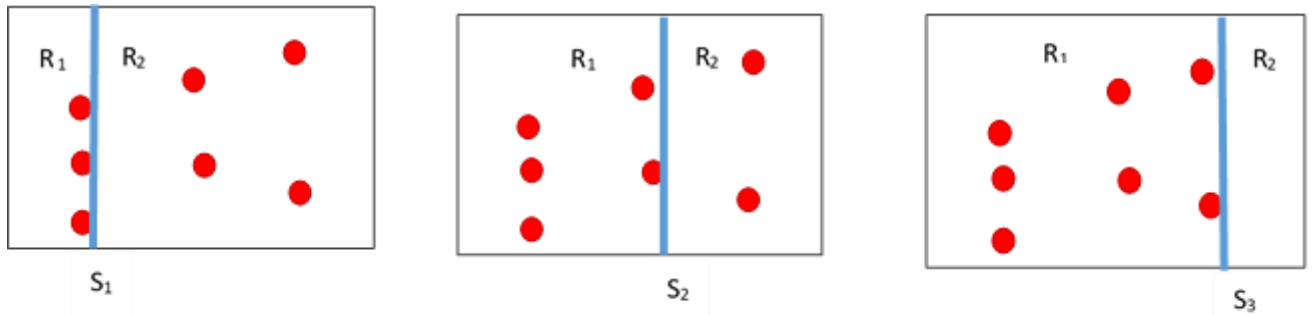
شکل ۶-۲: نمودار پراکنش داده‌های دو متغیره

شکل ۷-۲ و شکل ۸-۲ کلیه افرازهای ممکن و موثر را به ترتیب در راستای محورهای  $X_1$  و  $X_2$  نشان می‌دهند.

<sup>۱</sup>- Transmittal



شکل ۲-۷: چهار افراز ممکن در راستای متغیر  $x_2$



شکل ۲-۸: سه افراز ممکن در راستای متغیر  $x_1$

برای انتخاب بهترین افراز مرحله اول، به هر یک از هفت حالت فوق، یک اندازه نسبت داده

می‌شود. این اندازه به صورت رابطه زیر تعریف می‌گردد.



$$M = \sum_{R_1} (Y_i - \bar{Y})^2 + \sum_{R_2} (Y_i - \bar{Y})^2 \quad (3-2)$$

در این رابطه،  $Y_i$  متغیر پاسخ داده‌ها،  $\bar{Y}$  میانگین مقادیر متغیر پاسخ در هر قسمت و  $M$ ، برابر با مجموع توان‌های دوم خطا در نواحی  $R_1$  و  $R_2$  برای هریک از نمودارهای شکل ۲-۸ و شکل ۲-۷ است. افزایی که دارای کمترین مقدار  $M$  باشد، به عنوان بهترین افراز در مرحله اول انتخاب می‌شود و بدین ترتیب بهترین جهت و بهترین نقطه برای تفکیک فضا در گام اول بدست خواهد آمد. همچنین می‌توان نوشت

$$M = n_1 \text{Var}_1 + n_2 \text{Var}_2 \quad (4-2)$$

که در آن  $n_1 \text{Var}_1$  و  $n_2 \text{Var}_2$  به ترتیب واریانس مشاهدات واقع در نواحی  $R_1$  و  $R_2$  و همچنین  $n_1$  و  $n_2$  تعداد مشاهدات هریک از دو ناحیه افراز شده می‌باشد. بنابراین می‌توان  $M$  را مجموع واریانس‌های مشاهدات در دو ناحیه  $R_1$  و  $R_2$  دانست.

اکنون فضا به دو ناحیه تفکیک شده است. در مرحله دوم و مراحل بعدی فرآیند فوق در هریک از نواحی تولید شده تکرار می‌شود تا در نهایت فضای متغیرهای توضیحی به صورتی شبیه به قسمت قابل قبول شکل ۲-۳ (الف)، افراز شود.

### ۲-۳-۱-۳- اندازه درخت و هرس کردن<sup>۱</sup>

در این بخش به یک مطلب مهم درباره ساختار درخت رگرسیونی پرداخته می‌شود. تفکیک فضای متغیرها تا چه مرحله ای می‌تواند ادامه یابد یا به عبارت دیگر، هر درخت تا کجا می‌تواند رشد کند. اگر درخت به اندازه کافی رشد نکند ممکن است مدل مناسبی ارائه نشود و اگر هم بیش از حد

---

<sup>۱</sup>- Prune

رشد کند، احتمال این که بیش‌برآوردی<sup>۱</sup> رخ دهد، بسیار است. برای انتخاب اندازه درخت، نیاز به یک پارامتر در ساختار مدل است بطوری که این پارامتر مقدار بهینه اندازه درخت را بدست آورد.

در هر مرحله از رشد درخت، مدل حاصل از روش درخت رگرسیونی دقیق‌تر شده و مجموع توان‌های دوم خطا کاهش می‌یابد. راه توقف رشد درخت براساس کاهش مجموع توان‌های دوم خطا است، بدین ترتیب اگر در یک مرحله از رشد درخت، مجموع توان‌های دوم خطا دچار کاهش چندانی نشود، می‌توان از رشد درخت در آن مرحله صرف نظر کرد. این میزان کاهش در مورد مجموع توان‌های دوم خطا می‌تواند در اختیار کاربر باشد [۳۴]. در پایان این بخش باید یادآور شد که روش درخت رگرسیونی با بیش‌برآوردی همراه است، که برای جلوگیری از این مشکل، درخت نیاز به هرس دارد.

## ۲-۱۳-۴- معرفی روش جنگل‌های تصادفی

اساس روش RF وابسته به ماهیت روش درخت رگرسیونی است و از آنجایی که ساختار درخت رگرسیونی بطور کامل معرفی گردید، اینک می‌توان به معرفی روش RF پرداخت.

در روش RF مجموعه‌ای از درخت‌های رگرسیونی تشکیل می‌شوند و هر درخت مدلی را تولید می‌کند که مدل نهایی، برآیند یا ترکیبی از همه این مدل‌ها است. به عبارت دیگر هر یک از این درخت‌ها سهمی در مدل نهایی دارد. یکی از تفاوت‌های اساسی RF با درخت رگرسیونی آن است که در درخت رگرسیونی برای افزایش فضای متغیرها در هر مرحله، از کلیه متغیرها استفاده می‌گردد در حالی که در RF فقط از زیر مجموعه‌ای از متغیرها استفاده می‌شود. دیگر تفاوت اساسی این دو روش داده‌هایی است که، در ساخت مدل شرکت دارند. به این معنا که در درخت رگرسیونی، همه داده‌ها در ساخت مدل شرکت دارند، اما در روش RF تنها بخشی از داده‌ها در ساخت مدل استفاده می‌شوند. بطور کلی می‌توان گفت روش RF، ترکیبی از چندین درخت رگرسیونی است که در ساخت آن

---

<sup>۱</sup> -Over Estimation

چندین نمونه با جایگذاری از داده‌ها شرکت دارند و در هر درخت برای ساخت هر گره، تنها یک زیرمجموعه تصادفی از متغیرهای توضیحی شرکت می‌کنند. همچنین روش جنگل‌های تصادفی علاوه بر اینکه به عنوان یک روش مدل‌سازی به کار می‌رود به عنوان یک روش انتخاب متغیر نیز می‌تواند به کار برده شود.

## ۲-۱۳-۴-۱- الگوریتم روش جنگل‌های تصادفی

فرض کنید  $i = 1, \dots, N$  و  $(X_i, Y_i)$  مجموعه داده‌های مدل ساز (آموزش) باشند که در آن،  $X_i = (X_{i1}, \dots, X_{iM})$  برداری از  $M$  متغیر توضیحی و  $Y_i$  متغیر پاسخ متناظر آن است. اگر تعداد کل درخت‌های مدل با  $n_{tree}$  نشان داده شود، مراحل پنج‌گانه زیر بیانگر الگوریتم ساخت درخت  $i$ ام ( $i=1, \dots, n_{tree}$ ) است.

(۱) یک نمونه  $N$  تایی به روش با جایگذاری از مجموعه داده‌های آموزشی گرفته می‌شود. زیر مجموعه‌ای از داده‌های اصلی (حدود یک‌سوم از داده‌های آموزشی) که در این نمونه حضور ندارند را  $OOB^1$  نامیده که برای هر درخت نقش داده‌های آزمون را ایفا می‌کند. روش جداسازی داده‌های  $OOB$  روش نمونه‌گیری خودگردان است. با تکرار عملیات نمونه‌گیری، تعدادی مجموعه داده  $OOB$  از مجموعه آموزشی بوجود می‌آید که می‌توان برای هر مجموعه داده آموزشی یک درخت تصمیم محاسبه کرد.

(۲) به طور تصادفی از بین  $M$  متغیر توضیحی،  $m$  متغیر انتخاب می‌شود ( $m \ll M$ ). برای تقسیم فضای متغیرهای توضیحی به دو قسمت و براساس اصول درخت رگرسیونی، فقط از این  $m$  متغیر و نمونه  $N$  تایی انتخاب شده در گام ۱ استفاده می‌شود تا بهترین متغیر و بهترین نقطه افراز در اولین مرحله از مراحل بازگشتی بدست آید. محدودیت این مرحله آن است که تعداد مشاهدات

---

<sup>۱</sup> -Out of bag

موجود در هر یک از دو ناحیه باید بیشتر از  $n_r$  باشد و  $n_r$  پارامتری است که در اختیار کاربر بوده و بیانگر حداقل تعداد مشاهدات موجود در هر ناحیه است. معمولا در مدل رگرسیونی،  $m = \frac{M}{3}$  پیشنهاد می‌شود و در مدل رده‌بندی  $m = \sqrt{M}$  در نظر گرفته می‌شود.

(۳) برای هر یک از دو ناحیه تولید شده در گام ۲، مجدداً بطور تصادفی از بین  $M$  متغیر توضیحی،  $m$  متغیر انتخاب می‌گردد و با استفاده از همان نمونه  $N$  تایی انتخاب شده در گام ۱، روش درخت رگرسیون اعمال می‌گردد. این عمل منجر به افراز هر یک از نواحی موجود به دو قسمت می‌شود. لازم به ذکر است که افراز نواحی در صورتی انجام می‌شود که تعداد مشاهدات موجود در کل ناحیه بیشتر از  $2n_r$  باشد. یعنی اگر هر یک از نواحی دارای تعداد مشاهداتی کمتر از  $2n_r$  باشد نباید افراز دیگری روی هیچ یک از نواحی صورت گیرد.

(۴) گام ۳ برای تمام نواحی افراز شده تا زمانی تکرار می‌شود که تعداد مشاهدات در تمامی این نواحی کمتر از  $2n_r$  باشد.

(۵) حال درخت رگرسیونی  $\hat{m}$  طوری تشکیل شده است که فضای متغیرهای توضیحی به  $r_i$  ناحیه‌ی  $R_{1i}, R_{2i}, \dots, R_{r_i}$  تقسیم گردیده است و تعداد مشاهدات هر ناحیه کمتر از  $2n_r$  می‌باشد، مدل درخت رگرسیونی به دست آمده، به صورت زیر است.

$$\hat{f}_i(X) = \sum_{j=1}^{r_i} \hat{c}_j I_{R_{ij}}(X) \quad (5-2)$$

$$I_{R_{ij}}(X) = \begin{cases} 1, & X \in R_{ij} \\ 0, & X \notin R_{ij} \end{cases} \text{ و } \hat{c}_j = \hat{y}_j \text{ که در آن،}$$

از آن جایی که در روش RF تعداد  $n_{tree}$  درخت رگرسیونی وجود دارد، می‌توان گفت که تعداد  $n_{tree}$  مدل به صورت معادله (۵-۲) خواهیم داشت. اگر برای مقدار مشاهده شده‌ی  $x$  مدل خروجی درخت  $\hat{m}$  را به صورت  $\hat{f}_i(X)$  نشان دهیم، آن‌گاه برآورد متغیر پاسخ در این نقطه، با میانگین‌گیری از مقادیر  $\hat{f}_i(X)$  یعنی

$$\hat{y}(X) = \frac{1}{n_{tree}} \sum_{i=1}^{n_{tree}} \hat{f}_i(X) \quad (6-2)$$

به دست می آید.

لازم به ذکر است که در روش RF، تعداد درختها (پارامتر  $n_{tree}$ )، تعداد متغیرهای به تصادف انتخاب شده (پارامتر  $m$ ) و تعداد حداقل مشاهدات در هر افراز (پارامتر  $n_r$ )، پارامترهایی هستند که قابل تغییر بوده و توسط کاربر تعیین می گردند.

## ۲-۱۳-۴-۲- تعیین اهمیت متغیرها در روش جنگل های تصادفی

در ساختار روش RF امکانی وجود دارد که می توان میزان اهمیت متغیرها (VI)<sup>۱</sup> را در مدل، تعیین نموده و متغیرهایی که دارای نقش بیشتری در هر درخت و در مدل نهایی هستند شناسایی شوند. همانطور که در الگوریتم روش RF اشاره شد، برای تشکیل هر درخت، یک نمونه با جایگذاری از داده های اصلی مورد استفاده قرار می گیرد. همچنین داده هایی که در این نمونه حضور ندارند را OOB نامیدیم که به نوعی نقش داده های آزمایشی را برای ارزیابی آن درخت ایفا می کند. فرض کنید خطای پیش بینی  $Y$  برای داده های OOB در درخت  $i$ ام، با نماد  $EOOB_i$  نشان داده شود. برای تعیین اهمیت متغیر  $j$ ام، مقادیر این متغیر را بطور تصادفی  $n_{Perm}$  مرتبه ( $n_{Perm}$ ، پارامتری است که در اختیار کاربر می باشد) جا به جا کرده و مجددا خطا، به ازای مجموعه جدید محاسبه می شود که مقدار این خطا با نماد  $\widetilde{EOOB}_i^j$  نشان داده می شود. میزان اهمیت متغیر  $j$ ام در مدل درخت  $i$ ام با  $(VI_i(x^j))$  نشان داده و به صورت رابطه (۷-۲) تعریف می شود:

$$VI_i(x^j) = \widetilde{EOOB}_i^j - EOOB_i \quad (7-2)$$

سرانجام، میزان اهمیت متغیر  $j$ ام در مدل نهایی  $(VI(x^j))$  به صورت زیر است:

---

<sup>۱</sup>-Variable Importance

$$VI(x^j) = \frac{1}{ntree} \sum_{i=1}^{ntree} (VI_i(x^j)) \quad (۸-۲)$$

برای دقت بیشتر می‌توان جا به جا نمودن تصادفی مقادیر متغیر زام را به تعداد بیشتر از یک بار انجام داد و در هر بار خطای برآورد را محاسبه نمود و سرانجام مقدار  $\widetilde{EOOB}_1^j$  را به صورت میانگین خطاهای حاصله در نظر گرفت [۳۵ و ۳۲]. قابل ذکر است که اندازه اهمیت یک متغیر، به تنهایی قابل تفسیر نبوده و فقط برای رتبه بندی متغیرها براساس اهمیت آن‌ها در مدل به کار می‌رود.

## ۲-۱۳-۴-۳- مزیت‌های روش جنگل‌های تصادفی

- (۱) به دلیل شرکت کردن نمونه‌ای از داده‌ها در تشکیل مدل، این روش از لحاظ هزینه محاسبات مقرون به صرفه است. ضمن اینکه از دقت قابل قبولی نیز برخوردار می‌باشد.
- (۲) روش RF در مجموعه داده‌هایی که متغیرهای توضیحی زیادی وجود دارد، از کارایی بالایی برخوردار است. دلیل این ویژگی انتخاب تصادفی  $m$  متغیر از  $M$  متغیر کل در هر مرحله از رشد مدل می‌باشد.
- (۳) در این روش، با وجود این که هر درخت بطور کامل رشد می‌کند و هرس نمی‌گردد، با این حال مدل نهایی دچار بیش‌برآوردی نمی‌شود.
- (۴) این روش قابلیت استفاده در رده‌بندی داده‌های کیفی و پیش‌بینی داده‌های کمی را دارا می‌باشد.
- (۵) در صورت در اختیار داشتن تعداد مشاهدات کم، با توجه به نقش داده‌های OOB، می‌توان از آن به عنوان داده‌های آزمون استفاده کرد.
- (۶) این روش توانایی انتخاب متغیر را داراست و بر خلاف سایر روش‌های یادگیری ماشینی، نیازمند روشی دیگر برای انتخاب متغیر نیست [۳۲].

## ۲-۱۴- ارزیابی قدرت پیش‌بینی مدل

### ۲-۱۴-۱- استفاده از پارامترهای آماری

به منظور ارزیابی توانایی مدل‌ها در برآورد یا پیش‌بینی متغیر پاسخ مورد نظر، نیاز به معرفی بعضی معیارها است تا بتوان نتایج مدل‌های مختلف را با یکدیگر مقایسه کرد. در واقع این معیارها با اندازه‌گیری خطای برآورد یا دقت برآورد یک مدل، نتایج حاصل را قابل ارزیابی و قیاس‌پذیر با سایر نتایج می‌کنند. محققین در پایگاه‌های مختلف داده از معیارهای متفاوتی استفاده می‌کنند که به تعدادی از آن‌ها اشاره می‌شود.

**ضریب همبستگی<sup>۱</sup>:** ساده‌ترین راه برای بررسی میزان همبستگی دو یا چند متغیر، محاسبه

آماره ضریب همبستگی آن‌ها است. ضریب همبستگی دو متغیر  $x$  و  $y$  با رابطه (۲-۹) تعریف می‌شود:

$$R = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \quad (۲-۹)$$

مقدار این آماره بین ۱ تا -۱ متغیر است. مقدار بزرگتر قدر مطلق آن نشانه این است که ارتباط

خطی بیشتری میان متغیر وابسته و متغیرهای مستقل وجود دارد.

**ضریب تعیین<sup>۲</sup>:** ضریب تعیین یک معیار آماری است که بیانگر این است که چه درصدی از

تغییرات متغیر پاسخ را می‌توان توسط مدل ساخته شده توجیه کرد. ضریب تعیین را معمولاً با  $R^2$

نشان می‌دهند و مقدار آن بین صفر و یک است. در واقع  $R^2$  معیار نیست که بیانگر دقت مدل مورد نظر

در برآورد متغیر پاسخ می‌باشد. فرمول محاسباتی  $R^2$  به صورت زیر است.

---

<sup>۱</sup>-Correlation Coefficient

<sup>۲</sup>- Determination coefficient

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \quad (10-2)$$

که SSR<sup>۱</sup> در رابطه (۱۰-۲) بیانگر مجموع مربعات انحراف مقادیر پیش‌بینی شده‌ی متغیر وابسته از میانگین مقادیر آن است.

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (10-2)$$

SST<sup>۲</sup> در رابطه (۱۱-۲) نشانگر مجموع مربعات انحراف مقادیر واقعی متغیر وابسته از میانگین مقادیر آن است.

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (11-2)$$

SSE<sup>۳</sup> در رابطه (۱۲-۲) نیز مبین مجموع مربعات انحراف مقادیر واقعی متغیر وابسته از مقادیر پیش‌بینی شده برای آن است.

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (12-2)$$

که در این روابط  $\hat{y}_i$  مقدار پیش‌بینی شده‌ی متغیر وابسته،  $y_i$  مقدار واقعی متغیر وابسته و  $\bar{y}$  در هر رابطه، میانگین مقادیر متغیر وابسته است.

---

<sup>۱</sup> - Sum Square Regression

<sup>۲</sup> - Sum Square Total

<sup>۳</sup> - Sum Square Error



بنابراین با توجه به روابط فوق می توان نوشت:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (13-2)$$

طبق رابطه (۱۳-۲) اگر تمام مشاهدات بر روی خط برازش شده قرار گرفته باشند، یعنی به ازای تمام نقاط  $y_i = \hat{y}_i$  باشد، مقدار  $R^2$  برابر یک می شود و هرگونه انحرافی از این حالت باعث می شود که مقدار  $R^2$  از یک کوچکتر شود.

### خطای استاندارد پیش بینی<sup>۱</sup> (SEP)

SEP به صورت رابطه (۱۴-۲) تعریف می شود:

$$SEP = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}} \quad (14-2)$$

مجموع مربع باقیمانده ها<sup>۲</sup> (PRESS): برابر مجموع مربعات تفاوت بین مقدار

کمیت مشاهده شده ( $y_i$ ) و مقدار تخمین زده شده ( $\hat{y}_i$ ) است که مشابه رابطه (۱۲-۲) است.

$$PRESS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (15-2)$$

### میانگین مربع خطاها<sup>۳</sup> (MSE)

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n} \quad (16-2)$$

---

<sup>۱</sup> - Standard Error of Prediction

<sup>۲</sup> - Predictive Residual Sum of Squares

<sup>۳</sup> - Mean Square Error

میانگین خطای مطلق<sup>۱</sup> (MAE)

$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n} \quad (17-2)$$

میانگین خطای نسبی<sup>۲</sup> (MRE)

$$MRE = \frac{\sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|}{n} \times 100 \quad (18-2)$$

خطای نسبی پیش بینی<sup>۳</sup> (REP)

$$REP(\%) = \frac{100}{\bar{y}} \times \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}} \quad (19-2)$$

## ۲-۱۴-۲- استفاده از نمودار برگشتی

در نمودار برگشتی مقادیر پیش بینی شده‌ی کمیت مورد نظر برحسب مقادیر تجربی رسم می‌شود و به کمک ضریب تعیین ( $R^2$ ) به دست آمده از نمودار، پراکندگی نقاط اطراف خط برگشت تعیین می‌شود. محدوده‌ی تغییرات ضریب تعیین به دست آمده بین ۰ و ۱ است که هرچه قدر به ۱ نزدیک‌تر باشد، نتایج پیش‌بینی شده به مقادیر واقعی نزدیک‌ترند. اما اگر این مقدار برابر صفر باشد، بین داده‌ها هیچ گونه همبستگی وجود ندارد.

---

<sup>۱</sup> - Mean Absolute Error

<sup>۲</sup> - Mean Relative Error

<sup>۳</sup> - Relative Error of Prediction

## ۲-۱۴-۳- استفاده از نمودار خطای باقی مانده

منظور از عبارت خطای باقیمانده، اختلاف بین مقادیر پیش‌بینی شده و مقادیر تجربی است. اگر پراکندگی مقادیر در دو طرف خط صفر یکسان باشد، این امر نشان دهنده‌ی تصادفی بودن خطاست. ولی اگر عمده‌ی نقاط در این نمودار، در یک طرف خط صفر باشند، این بدان معناست که خطای جهت‌داری رخ داده است.

## ۲-۱۴-۴- استفاده از سری تست

اعتبار و اهمیت مدل‌های پیش‌بینی کننده وقتی مشخص می‌شود که بتوانند خاصیت مورد نظر را برای مولکول‌های جدیدی که در دسته‌ی آموزش موجود نیستند پیش‌بینی کنند. بدین منظور از ابتدای کار تعدادی از مولکول‌ها به عنوان سری تست خارجی کنار گذاشته می‌شود و در نهایت مدل منتخب بهینه شده، جهت پیش‌بینی خاصیت مورد نظر مولکول‌های موجود در سری تست که در مدل‌سازی استفاده نشده‌اند، به کار گرفته می‌شود.

## ۲-۱۴-۵- استفاده از روش رد مرحله‌ای تک تک

در این روش هم برای تعیین اعتبار مدل‌های پیش‌بینی کننده هر بار یک ترکیب به عنوان سری آزمون کنار گذاشته می‌شود و در شرایط بهینه به دست آمده برای سری ارزیابی یا سری آموزش با ترکیبات باقی‌مانده مدل‌سازی صورت می‌گیرد، سپس مدل بدست آمده برای پیش‌بینی فعالیت ترکیب کنار گذاشته شده به کار گرفته می‌شود و این فرآیند برای تمام اعضای سری داده‌ها تکرار می‌گردد.



۳- فصل سوم

مطالعه ارتباط کمی ساختار-فعالیت

مشتقات پیریدینوز

## ۳-۱- مدل سازی فعالیت دارویی مشتقات پیریدینون

در این فصل، روش جنگل تصادفی و شبکه عصبی مصنوعی به منظور مطالعه ارتباط کمی ساختار-فعالیت برای پیش‌بینی فعالیت ضد HIV مشتقات ۴-(۱و۵-دی بنزیل-۱و۲-دی هیدرو-۲-اکسو پیریدین-۳-یل)-۲-هیدروکسی-۴-اکسوبوت-۲-انوئیک اسید<sup>۱</sup> به کار گرفته شده است. و توانایی مدل حاصله در پیش‌بینی فعالیت‌های ضد ایدز این ترکیبات مورد ارزیابی قرار گرفته است. بطور کلی این تحقیق شامل معرفی سری داده‌ها، بهینه سازی ساختار مولکول‌ها و محاسبه توصیف کننده‌های مولکولی با استفاده از روش سهم گروه و رگرسیون مرحله‌ای و همچنین مدل سازی توسط روش‌های جنگل‌های تصادفی و شبکه عصبی مصنوعی و ارزیابی مدل‌های برتر است.

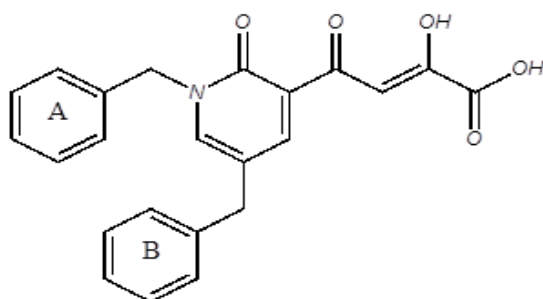
## ۳-۲- سری داده‌ها

سری داده‌ها از نتایج تجربی ۵۳ ترکیب فعال از مشتقات پیریدینون می‌باشد. اسکلت اصلی این ترکیبات در شکل ۳-۱ نشان داده شده است. جدول ۳-۱ استخلاف‌های مختلف تمام این ترکیبات را نشان می‌دهد. همچنین مقادیر  $pEC_{50}$  (منهای لگاریتم  $EC_{50}$ <sup>۲</sup> بر حسب مولار) این ترکیبات به عنوان معیاری کمی از فعالیت ضد HIV این ترکیبات در جدول ۳-۱ گزارش شده است که از مرجع [۳۶] استخراج گردیده است.  $EC_{50}$  غلظتی از دارو می‌باشد که تا ۵۰٪ تکثیر ویروس HIV-1 را مهار می‌کند.

---

<sup>۱</sup> - 4-(1,5-dibenzyl-1,2-dihydro-2-oxopyridin-3-yl)-2-hydroxy-4-oxobut-2-enoic acid

<sup>۲</sup> -Effective Concentration



شکل ۱-۳: اسکلت اصلی ترکیب مورد مطالعه

جدول ۱-۳: استخلاف‌های ترکیب مورد مطالعه

شماره ترکیب	حلقه آرلی A			حلقه آرلی B			pEC <sub>50</sub>
	ارتو	متا	پارا	ارتو	متا	پارا	
۱	H	H	H	H	H	H	۵/۶۷
۲	F	H	H	H	H	H	۶/۰۴
۳	H	H	F	H	H	H	۶/۰۰
۴	F	H	F	H	H	H	۶/۲۲
۵	F	H	H	F	H	H	۶/۱۵
۶	H	H	OMe	H	H	H	۵/۷۴
۷	H	H	F	H	H	F	۶/۰۹
۸	H	H	Me	H	H	F	۵/۷۹
۹	H	Cl	F	H	H	H	۶/۰۹
۱۰	F	H	F	F	H	F	۶/۵۲
۱۱	H	F	H	H	H	H	۵/۶۹
۱۲	H	Cl	H	H	Cl	H	۵/۹۵
۱۳	H	Cl	F	H	Cl	H	۵/۹۲
۱۴	F	H	H	H	H	F	۶/۱۵
۱۵	H	Cl	F	H	Cl	F	۵/۶۱
۱۶	H	H	Me	H	Cl	H	۵/۶۰
۱۷	H	H	Me	F	H	H	۶/۱۵
۱۸	F	Cl	F	H	H	H	۶/۲۲
۱۹	H	H	F	H	Cl	H	۵/۸۵
۲۰	2,6-di-F	H	H	H	H	F	۵/۸۵
۲۱	H	H	F	F	H	H	۶/۰۹

شماره ترکیب	حلقه آريلي A			حلقه آريلي B			pEC <sub>50</sub>
	ارتو	متا	پارا	ارتو	متا	پارا	
۲۲	F	Cl	H	F	Cl	H	۵/۶۱
۲۳	H	H	Cl	H	Cl	H	۵/۶۰
۲۴	H	Cl	H	H	Cl	H	۵/۸۵
۲۵	F	Cl	H	H	Me	H	۵/۸۵
۲۶	F	H	F	H	F	H	۶/۲۲
۲۷	H	Me	H	H	Cl	F	۵/۶۵
۲۸	F	Cl	H	H	H	H	۶/۱۵
۲۹	F	5-Cl	H	H	H	H	۶/۰.۴
۳۰	F	Cl	H	H	Cl	H	۵/۸۵
۳۱	F	H	F	H	Cl	F	۵/۸۵
۳۲	H	Cl	H	H	H	F	۵/۹۲
۳۳	H	Cl	H	F	H	H	۶/۱۵
۳۴	2,6-di-F	H	H	H	Cl	F	۵/۲۹
۳۵	2,6-di-F	H	H	H	H	H	۵/۸۲
۳۶	H	H	H	F	H	F	۶/۰.۹
۳۷	F	H	H	H	Cl	F	۵/۶۵
۳۸	H	F	H	F	H	H	۵/۷۲
۳۹	F	H	H	H	Cl	H	۵/۷۲
۴۰	F	F	H	H	H	H	۶/۰.۴
۴۱	H	Me	H	F	H	H	۶/۰.۴
۴۲	F	F	H	H	H	F	۶/۱۵
۴۳	H	H	CN	F	H	H	۵/۹۵
۴۴	F	F	H	F	H	H	۶/۰.۹
۴۵	H	F	H	H	H	F	۶/۰.۰
۴۶	F	H	F	Cl	H	H	۶/۰.۹
۴۷	H	Cl	H	F	Cl	H	۵/۸۲
۴۸	H	H	F	F	H	F	۶/۲۲
۴۹	H	Me	H	F	H	F	۶/۱۵
۵۰	H	H	H	Me	H	F	۵/۸۸
۵۱	H	H	F	Me	H	F	۶/۰.۰
۵۲	F	H	H	F	H	F	۶/۳۰
۵۳	F	H	H	Me	H	F	۵/۸۸



## ۳-۳- نرم‌افزارهای مورد استفاده

دانش کمومتریکس بسته‌های نرم‌افزاری متنوعی را برای انجام تمام مراحل مدل‌سازی به کار گرفته است که در ادامه به اختصار بسته‌های نرم‌افزاری استفاده شده در این تحقیق معرفی می‌شوند.

### ۳-۳-۱- بسته نرم‌افزاری Hyperchem

از بسته نرم‌افزاری Hyperchem [۳۷] برای رسم شکل مولکول‌ها و بهینه‌سازی ساختار با استفاده از روش‌های کوانتومی و مکانیکی، استفاده می‌شود. به کمک این برنامه می‌توان طول پیوند، زاویه پیوندی و زوایای پیچشی را در مولکول تعیین کرد. داده‌های حاصل از این نرم‌افزار را می‌توان به عنوان ورودی به سایر نرم‌افزارها معرفی نمود. به عنوان مثال ساختار سه بعدی مولکول‌ها که با استفاده از روش مکانیک کوانتومی در Hyperchem بهینه شده است، را می‌توان جهت محاسبه توصیف‌کننده‌های مختلف به وسیله نرم‌افزار دراگون<sup>۱</sup>، به کار برد. همچنین به کمک این نرم‌افزار می‌توان تعدادی از توصیف‌کننده‌ها از جمله حجم مولی و قطبش‌پذیری را محاسبه کرد.

### ۳-۳-۲- بسته نرم‌افزاری دراگون

نرم‌افزار Dragon5.5 امکان محاسبه ۳۲۲۴ توصیف‌کننده مختلف را فراهم می‌کند. جهت محاسبه توصیف‌کننده به کمک این نرم‌افزار لازم است ساختار هندسی بهینه مولکول مورد استفاده قرار گیرد. برای این منظور می‌توان ساختار بهینه مولکول‌ها را به صورت فایل‌هایی با فرمت mol, self و hin به عنوان اطلاعات ورودی به کار برد.

---

<sup>۱</sup>- Dragon

### ۳-۳-۳- بسته نرم‌افزاری SPSS<sup>۱</sup>

SPSS [۳۸] که نخستین نسخه آن در سال ۱۹۷۰ توسط جمعی از فارغ التحصیلان دانشگاه استانفورد آمریکا ارائه شده، امکان تجزیه و تحلیل آماری داده‌ها را فراهم می‌آورد. برخی از قابلیت‌های این بسته نرم‌افزاری عبارتند از:

- (۱) تعیین تعداد فراوانی‌های هر یک از گروه‌ها در یک متغیر
- (۲) محاسبه میانگین ساده برای داده‌ها
- (۳) نمایش اطلاعات به صورت متنوع در قالب نمودار و جدول
- (۴) انجام رگرسیون تک متغیره و چند متغیره

### ۳-۳-۴- نرم‌افزار MATLAB

نرم‌افزار MATLAB [۳۹] یکی از جامع‌ترین و کارآمدترین نرم‌افزارهای علمی و محاسباتی است که طی چند سال گذشته تهیه و به بازار عرضه شد و در سال‌های اخیر با تدوین نسخه‌های جدیدتر و کامل‌تر روز به روز بر غنای آن افزوده می‌شود.

MATLAB<sup>۲</sup> به معنای آزمایشگاه ماتریس است. ورودی‌ها اساساً به صورت ماتریس در نظر گرفته می‌شوند و هیچ نیازی به مشخص کردن ابعاد ماتریس نمی‌باشد. در MATLAB حتی اعداد اسکالر، ماتریس‌های (x) به حساب می‌آیند و بردارها، حالت خاصی از ماتریس‌های سطری یا ستونی در نظر گرفته می‌شوند. از جمله کاربردهای جالب آن که مورد توجه شیمی‌دانان قرار گرفته است شبکه‌های عصبی مصنوعی، جنگل‌های تصادفی و ماشین بردار پشتیبان است که در این پژوهش از آن‌ها استفاده شده است. اطلاعات شیمیایی به عنوان ورودی و به صورت یک ماتریس به نرم‌افزار داده

---

<sup>۱</sup> - Statistical Package for the Social Science

<sup>۲</sup> - Matrix Laboratory

می‌شود. در محیط این برنامه با استفاده از آرایه‌ها و فرامان‌های موجود امکان مدل‌سازی غیرخطی فعالیت بیولوژیکی با ساختار ترکیبات فراهم می‌گردد.

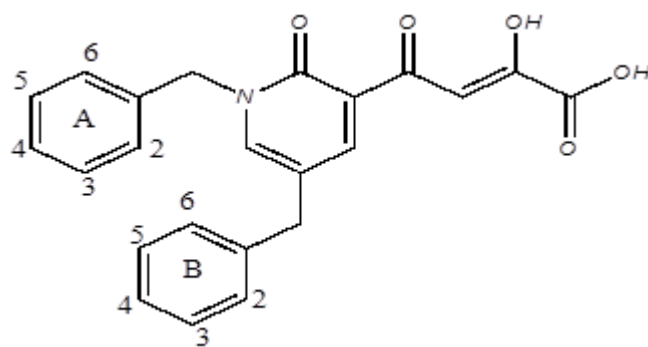
### ۳-۴- بهینه‌سازی ساختمان هندسی مولکول‌ها و محاسبه توصیف‌کننده‌ها

زمانی که متغیرها را به روش‌هایی مانند رگرسیون مرحله‌ای انتخاب می‌کنیم نیاز است که ساختار مولکول‌ها را بهینه کنیم. ساختار مولکول‌ها با استفاده از نرم افزار Hyperchem رسم و با روش نیمه تجربی AM1 برای رسیدن جذر میانگین مربعات گرادیان انرژی به  $0.001$  کیلوکالری بر مول بهینه شدند. سپس بوسیله نرم افزار دراگون، تعداد ۳۲۲۴ توصیف‌کننده با استفاده از این ساختارهای بهینه شده محاسبه شد.

### ۳-۵- انتخاب بهترین توصیف‌کننده‌ها

#### ۳-۵-۱- معرفی توصیف‌کننده‌ها به روش سهم گروه (GC)

در این بخش برای انتخاب توصیف‌کننده‌ها از روش سهم گروه استفاده شده است. برای این منظور، مولکول‌ها به اجزای مختلف شکسته شدند که در نهایت برای ۵۳ مولکول، گروه‌های ساختاری به صورتی که در جدول ۳-۳ نشان داده خواهد شد، بدست آمد. توصیف‌کننده‌های موجود در این روش بیانگر نوع استخلاف و جایگاه موقعیت‌های استخلاف شده در حلقه‌های آروماتیک A و B می‌باشد. برای نشان دادن نوع موقعیت‌هایی که دارای استخلاف هستند، این جایگاه‌ها همانطور که در شکل ۳-۲ نشان داده شده در حلقه‌های آروماتیک A و B شماره گذاری گردید. برای مثال استخلاف در جایگاه شماره دو در حلقه آریلی A با نماد  $A_2$  نمایش داده شد به همین شیوه  $A_3, A_4, A_5, A_6$  و  $B_2, B_3, B_4, B_5, B_6$  به دست آمدند.



شکل ۲-۳: نمایش نوع موقعیت‌ها (جایگاه‌های) استخلاف شده

همچنین برای تشخیص اینکه توصیف کننده دارای چه نوعی استخلافی است مطابق جدول ۲-۳ از اعداد استفاده شد. لازم به ذکر است که برای توصیف کننده‌های دارای استخلاف گیرنده الکترون علامت مثبت و برای توصیف کننده‌هایی با استخلاف دهنده الکترون علامت منفی قرار داد شد. علاوه بر این، مقدار اعداد نسبت داده شده به توصیف کننده‌ها براساس قدرت الکترون کشندگی است و به‌طور قراردادی نسبت داده شده است، بطور مثال برای فلئور که قویترین کشنده الکترون است عدد ۱ و برای کلر که قدرت الکترون کشندگی آن کمتر از فلئور است عدد ۰/۵ در نظر گرفته شد.

جدول ۲-۳: مقادیر نسبت داده شده به توصیف کننده‌ها

نوع توصیف کننده	مقدار توصیف کننده
<b>F</b>	۱
<b>CN</b>	۰/۷
<b>Cl</b>	۰/۵
<b>H</b>	۰
<b>Me</b>	-۰/۵
<b>OMe</b>	-۱

جدول ۳-۳: توصیف‌کننده‌های سهم گروه

Number	A2	A3	A4	A5	A6	B2	B3	B4	B5	B6
۱	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
۲	۱	۰	۰	۰	۰	۰	۰	۰	۰	۰
۳	۰	۰	۱	۰	۰	۰	۰	۰	۰	۰
۴	۱	۰	۱	۰	۰	۰	۰	۰	۰	۰
۵	۱	۰	۰	۰	۰	۰	۰	۰	۰	۱
۶	۰	۰	-۱	۰	۰	۰	۰	۰	۰	۰
۷	۰	۰	۱	۰	۰	۰	۰	۱	۰	۰
۸	۰	۰	-۰/۵	۰	۰	۰	۰	۱	۰	۰
۹	۰	۰/۵	۱	۰	۰	۰	۰	۰	۰	۰
۱۰	۱	۰	۱	۰	۰	۰	۰	۱	۰	۱
۱۱	۰	۱	۰	۰	۰	۰	۰	۰	۰	۰
۱۲	۰	۰/۵	۰	۰	۰	۰	۰	۰	۰/۵	۰
۱۳	۰	۰/۵	۰	۰	۰	۰	۰	۱	۰/۵	۰
۱۴	۱	۰	۰	۰	۰	۰	۰	۱	۰	۰
۱۵	۰	۰/۵	۱	۰	۰	۰	۰	۱	۰/۵	۰
۱۶	۰	۰	-۰/۵	۰	۰	۰	۰	۰	۰/۵	۰
۱۷	۰	۰	-۰/۵	۰	۰	۰	۰	۰	۰	۱
۱۸	۱	۰/۵	۱	۰	۰	۰	۰	۰	۰	۰
۱۹	۰	۰	۱	۰	۰	۰	۰	۰	۰/۵	۰
۲۰	۱	۰	۰	۰	۱	۰	۰	۱	۰	۰
۲۱	۰	۰	۱	۰	۰	۰	۰	۰	۰	۱
۲۲	۱	۰/۵	۰	۰	۰	۰	۰	۰	۰/۵	۱
۲۳	۰	۰	۰/۵	۰	۰	۰	۰	۰	۰/۵	۰
۲۴	۰	۰/۵	۰	۰	۰	۰	۰	۰	۰/۵	۰
۲۵	۱	۰/۵	۰	۰	۰	۰	۰	۰	-۰/۵	۰
۲۶	۱	۰	۱	۰	۰	۰	۰	۰	۰	۱
۲۷	۰	-۰/۵	۰	۰	۰	۰	۰	۱	۰/۵	۰
۲۸	۱	۰/۵	۰	۰	۰	۰	۰	۰	۰	۰
۲۹	۱	۰	۰	۰/۵	۰	۰	۰	۰	۰	۰
۳۰	۱	۰/۵	۰	۰	۰	۰	۰	۰	۰/۵	۰
۳۱	۱	۰	۱	۰	۰	۰	۰	۱	۰/۵	۰

ادامه جدول (۳-۳)

Number	A2	A3	A4	A5	A6	B2	B3	B4	B5	B6
۳۲	۰	۰/۵	۰	۰	۰	۰	۰	۱	۰	۰
۳۳	۰	۰/۵	۰	۰	۰	۰	۰	۰	۰	۱
۳۴	۱	۰	۰	۰	۱	۰	۰	۱	۰/۵	۰
۳۵	۱	۰	۰	۰	۱	۰	۰	۰	۰	۰
۳۶	۰	۰	۰	۰	۰	۰	۰	۱	۰	۱
۳۷	۱	۰	۰	۰	۰	۰	۰	۱	۰/۵	۰
۳۸	۰	۱	۰	۰	۰	۰	۰	۰	۰	۱
۳۹	۱	۰	۰	۰	۰	۰	۰	۰	۰/۵	۰
۴۰	۱	۱	۰	۰	۰	۰	۰	۰	۰	۰
۴۱	۰	-۰/۵	۰	۰	۰	۰	۰	۰	۰	۱
۴۲	۱	۱	۰	۰	۰	۰	۰	۱	۰	۰
۴۳	۰	۰	۰/۷	۰	۰	۰	۰	۰	۰	۱
۴۴	۱	۱	۰	۰	۰	۰	۰	۰	۰	۱
۴۵	۰	۱	۰	۰	۰	۰	۰	۱	۰	۰
۴۶	۱	۰	۱	۰	۰	۰	۰	۰	۰	۰/۵
۴۷	۰	۰/۵	۰	۰	۰	۰	۰	۰	۰/۵	۱
۴۸	۰	۰	۱	۰	۰	۰	۰	۱	۰	۱
۴۹	۰	-۰/۵	۰	۰	۰	۰	۰	۱	۰	۱
۵۰	۰	۰	۰	۰	۰	۰	۰	۱	۰	-۰/۵
۵۱	۰	۰	۱	۰	۰	۰	۰	۱	۰	-۰/۵
۵۲	۱	۰	۰	۰	۰	۰	۰	۱	۰	۱
۵۳	۱	۰	۰	۰	۰	۰	۰	۱	۰	۰/۵

همانطور که در جدول ۳-۳ مشاهده می‌شود همه توصیف کننده‌های ستون‌های B<sub>2</sub> و B<sub>3</sub> صفر هستند، لذا پس از حذف داده‌های ستون‌های B<sub>2</sub> و B<sub>3</sub> ، داده‌های ۸ ستون باقی‌مانده به عنوان ورودی برای روش جنگل تصادفی استفاده شد.

### ۳-۵-۲- توصیف‌کننده‌های حاصل از دراگون و انتخاب بهترین توصیف‌کننده‌ها با

#### روش رگرسیون مرحله‌ای (SR)

با توجه به این که تعداد زیاد توصیف‌کننده باعث پیچیدگی محاسبات شده و تعدادی از توصیف‌کننده‌ها حاوی اطلاعات یکسانی هستند، تعداد توصیف‌کننده‌ها باید به نوعی کاهش یابد. به همین منظور در نرم افزار دراگون هنگام ذخیره سازی توصیف‌کننده‌ها، آن‌ها با همبستگی  $0/96$  ذخیره گردید که با این کار توصیف‌کننده‌هایی که برای تمام مولکول‌ها مقادیر ثابت یا تقریباً ثابت داشتند، حذف شدند و همچنین یکی از دو توصیف‌کننده‌ای که همبستگی بزرگتر از  $0/96$  داشت حذف گردید. در نهایت از  $3224$  توصیف‌کننده،  $600$  توصیف‌کننده باقی ماند. با به کارگیری نرم افزار SPSS و با اجرای رگرسیون مرحله‌ای برای  $600$  توصیف‌کننده‌های یاد شده به عنوان متغیرهای مستقل، و فعالیت به عنوان متغیر وابسته، تعداد  $9$  توصیف‌کننده انتخاب شدند که نام و طبقه آن‌ها در جدول ۳-۴ نشان داده شده است. همچنین جدول ۳-۵ ماتریس همبستگی بین این توصیف‌کننده‌ها را ارائه می‌دهد، نتایج این جدول نشان می‌دهد که بین این توصیف‌کننده‌ها همبستگی معناداری وجود ندارد.

جدول ۳-۴: توصیف کننده‌های انتخاب شده توسط رگرسیون مرحله‌ای

No	Symbol	Class	Meaning
۱	B07NCI	2D binary fingerprints	Presence/absence of N-Cl at topological distance 07
۲	Mor25m	3D-MoRSE descriptors	3D-MoRSE- signal 25/Weighted by atomic masses
۳	MATS8m	2D autocorrelation	Moaran autocorrelation -lag8//Weighted by atomic masses
۴	Mor29v	3D-MoRSE descriptors	3D-MoRSE- signal 29/Weighted by atomic van der waals volumes
۵	RDF110u	RDF descriptors	Radial Distribution Function - 11.0/unweighted
۶	TF..Cl	Topological descriptors	Sum of topological distances between F..Cl
۷	G1p	WHIM descriptors	1st component symmetry directional WHIM index/ Weighted by atomic polarizabilities
۸	CIC5	Information indices	Complementary information content(neighborhood symmetry of 5-order)
۹	RDF095u	RDF descriptors	Radial Distribution Function -9.5/unweighted

جدول ۳-۵: ماتریس همبستگی کل توصیف کننده‌های انتخاب شده توسط رگرسیون مرحله‌ای

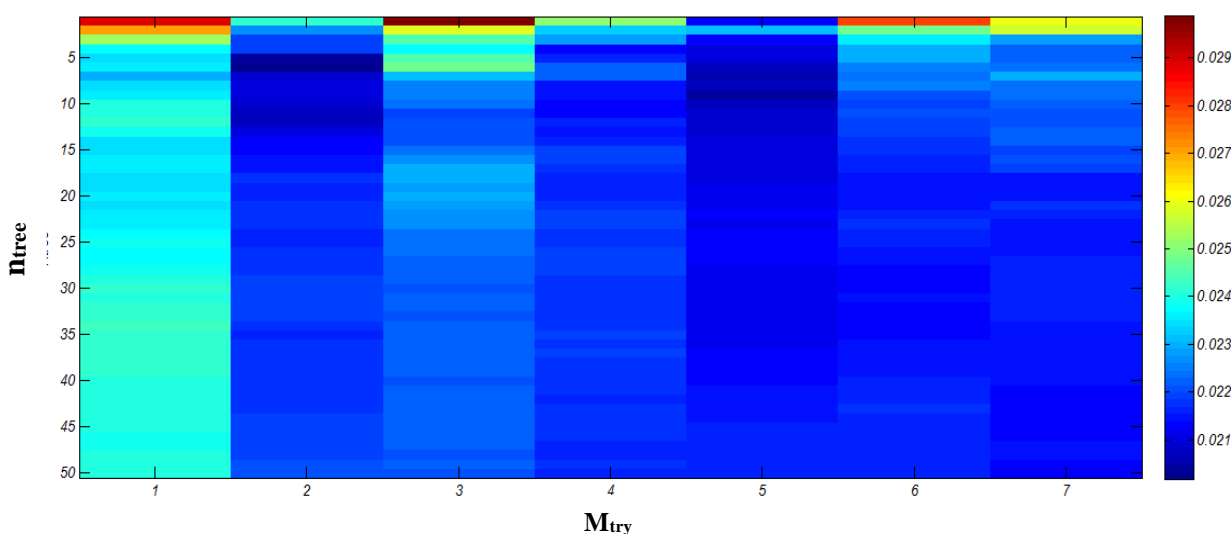
	B07NCI	Mor25m	MATS8m	Mor29v	RDF110u	TF..Cl	G1p	CIC5	RDF095u
B07NCI	۱								
Mor25m	۰/۵۳۴	۱							
MATS8m	-۰/۳۹۹	-۰/۱۰۷	۱						
Mor29v	-۰/۱۸۳	۰/۰۶۷	۰/۱۹۹	۱					
RDF110u	۰/۰۴۵	۰/۰۵۲	۰/۰۴۱	-۰/۰۰۱	۱				
TF..Cl	۰/۶۰۱	۰/۵۳۳	-۰/۲۶۸	-۰/۰۳۷	-۰/۱۰۵	۱			
G1p	-۰/۰۵۸	-۰/۰۲۵	-۰/۰۳۱	۰/۱۶۳	-۰/۲۷۷	-۰/۰۹۹	۱		
CIC5	-۰/۳۸۵	-۰/۵۲۱	-۰/۱۷۴	-۰/۱۱۷	۰/۱۶۶	-۰/۴۱۲	-۰/۱۱۹	۱	
RDF095u	۰/۰۱۳	۰/۰۷۲	۰/۰۳۸	-۰/۰۷۰	۰/۲۳۰	-۰/۰۷۱	-۰/۰۹۳	-۰/۰۰۹	۱



### ۳-۶- مدل سازی جنگل های تصادفی

۳-۶-۱- مدل سازی جنگل های تصادفی با توصیف کننده های سهم گروه (GC-RF)

ابتدا کل داده ها به دو سری آموزش (۴۳ ترکیب) و سری آزمون (۱۰ ترکیب) تقسیم شدند. برای ایجاد مدل جنگل های تصادفی از داده های سری آموزش استفاده گردید. همانطور که در بخش (۲-۱۲-۴-۱) گفته شد،  $\frac{1}{3}$  داده ها در ساخت درخت شرکت ندارند، این داده ها OOB هستند که برای بهینه سازی پارامترهای موثر در نظر گرفته شدند. برای دستیابی به بهترین مدل ابتدا سه پارامتر مهم و تاثیر گذار یعنی تعداد درختان ( $n_{tree}$ )، تعداد توصیف کننده های انتخاب شده در هر مرحله افزایش ( $M_{try}$ ) و تعداد مشاهدات باقی مانده در هر گره (اندازه گره) بهینه گردیدند. برای بهینه نمودن این پارامترها، تعداد درختها از ۱۰۰ تا ۵۰۰ با گام صد، تعداد توصیف کننده های انتخاب شده در هر مرحله افزایش ( $M_{try}$ ) از ۲ تا ۸ با گام یک و تعداد مشاهدات باقی مانده در هر گره (Node Size) از ۱ تا ۱۰ با گام یک، بطور همزمان تغییر داده شد و در هر مرحله مقدار MSE مربوط به مجموعه OOB محاسبه گردید. نتایج بهینه سازی پارامترهای  $n_{tree}$  و  $M_{try}$  در شکل ۳-۳ نشان داده شده است همچنین جدول ۳-۶ چندین حالت که کمترین MSE را برای OOB دارد نشان می دهد.



شکل ۳-۳: بهینه سازی پارامترهای  $n_{tree}$  و  $M_{try}$  در Node size بهینه

همانطور که در شکل ۳-۳ مشاهده می‌شود نقاط آبی رنگ خطای کمتری را نشان می‌دهد.

جدول ۳-۶: کمترین مقادیر MSE همراه با  $M_{try}$  و  $n_{tree}$  متناظر با آنها

$n_{tree}$	$M_{try}$	Node Size	MSE OOB
۲۰۰	۸	۱	۰/۰۰۲۴۱
۱۰۰	۲	۲	۰/۰۰۳۵۵
۱۰۰	۸	۳	۰/۰۰۴۹۸
۱۰۰	۸	۴	۰/۰۰۶۷۰
۱۰۰	۸	۵	۰/۰۰۸۲۰
۱۰۰	۸	۶	۰/۰۰۹۶۵
۲۰۰	۸	۷	۰/۰۱۰۴۶
۱۰۰	۸	۸	۰/۰۱۱۸۱
۱۰۰	۸	۹	۰/۰۱۲۵۵
۱۰۰	۸	۱۰	۰/۰۱۳۴۷

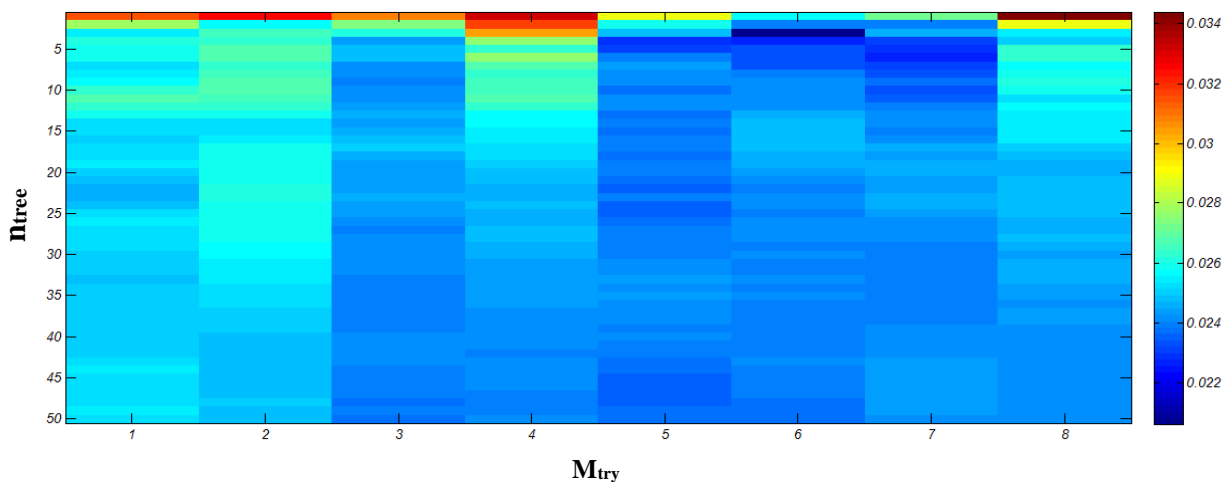
با توجه به نتایج جدول، تعداد درختان = ۲۰۰ و تعداد توصیف کننده‌های انتخاب شده در هر مرحله افراز = ۸ و تعداد گره = ۱ به عنوان مقادیر بهینه انتخاب شدند.

### ۳-۶-۲- مدل‌سازی جنگل‌های تصادفی با توصیف کننده‌های انتخاب

#### شده به وسیله روش رگرسیون مرحله‌ای (SR)

ابتدا کل داده‌ها به دو سری آموزش (۴۳ ترکیب) و سری آزمون (۱۰ ترکیب) تقسیم شدند. برای ایجاد مدل جنگل‌های تصادفی از داده‌های سری آموزش استفاده گردید،  $\frac{1}{3}$  داده‌ها در ساخت درخت شرکت ندارند، این داده‌ها OOB هستند که برای بهینه سازی پارامترهای موثر در نظر گرفته شدند. برای دستیابی به بهترین مدل ابتدا سه پارامتر مهم و تاثیر گذار یعنی تعداد درختان ( $n_{tree}$ )، تعداد توصیف کننده‌های انتخاب شده در هر مرحله افراز ( $M_{try}$ ) و تعداد مشاهدات باقی مانده در هر گره (Node Size) برای داده‌های OOB بهینه گردیدند. برای بهینه نمودن این پارامترها، تعداد درخت

از ۱۰۰ تا ۵۰۰ با گام صد، تعداد توصیف کننده‌های انتخاب شده در هر مرحله افراز ( $M_{try}$ ) از ۲ تا ۹ با گام یک و تعداد مشاهدات باقی مانده در هر گره (Node Size) از ۱ تا ۱۰ با گام یک، بطور همزمان تغییر داده شد و در هر مرحله مقدار MSE مربوط به مجموعه OOB محاسبه گردید. نتایج بهینه سازی پارامترهای  $M_{try}$  و  $n_{tree}$  در شکل ۳-۴ نشان داده شده است که برای نقاط آبی مقدار mse کمتر است و جدول ۳-۷ چندین حالت که کمترین MSE را برای OOB دارد نشان می‌دهد.



شکل ۳-۴: بهینه سازی پارامترهای  $M_{try}$  و  $n_{tree}$  در Node size بهینه

جدول ۳-۷: کمترین مقادیر MSE همراه با  $M_{try}$  و  $n_{tree}$  متناظر با آنها

$n_{tree}$	$M_{try}$	Node Size	MSE OOB
۳۰۰	۴	۱	۰/۰۰۲۹۹
۱۰۰	۵	۲	۰/۰۰۳۱۵
۱۰۰	۷	۳	۰/۰۰۳۵۴
۲۰۰	۹	۴	۰/۰۰۴۰۰
۳۰۰	۹	۵	۰/۰۰۵۰۹
۳۰۰	۹	۶	۰/۰۰۶۰۱
۵۰۰	۸	۷	۰/۰۰۷۰۶
۲۰۰	۹	۸	۰/۰۰۷۸۴
۴۰۰	۹	۹	۰/۰۰۹۱۲
۳۰۰	۸	۱۰	۰/۰۱۰۰۷

با توجه به نتایج به دست آمده، تعداد درختان = ۳۰۰ و تعداد توصیف کننده‌های انتخاب شده در هر مرحله افراز = ۴ و تعداد گره = ۱ به عنوان مقادیر بهینه انتخاب شدند.

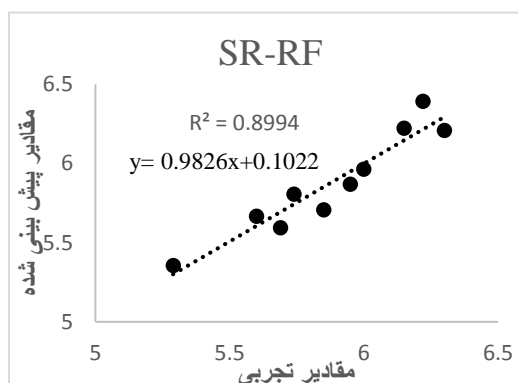
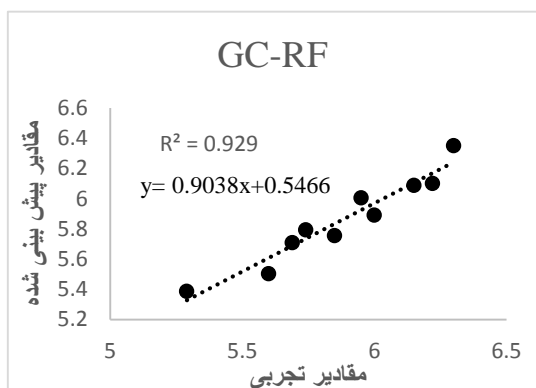
### ۳-۶-۳- ارزیابی مدل جنگل‌های تصادفی

#### ۳-۶-۳-۱- ارزیابی مدل‌های RF با استفاده از داده‌های سری آزمون

اهمیت مدل‌های پیش‌بینی وقتی مشخص می‌گردد که مدل حاصله بتواند فعالیت مولکول‌هایی را که در مدل‌سازی به کار نرفته‌اند پیش‌بینی کند. بدین منظور مدل‌های منتخب، برای پیش‌بینی فعالیت داده‌های سری آزمون به کار برده شدند. جدول ۳-۸ نتایج پیش‌بینی حاصل از ارزیابی مدل‌های جنگل تصادفی (GC-RF و SR-RF) با استفاده از داده‌های سری آزمون و شکل ۳-۵ نیز نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی را برای داده‌های سری آزمون نشان می‌دهد.

جدول ۳-۸: نتایج حاصل از ارزیابی مدل‌های جنگل تصادفی با استفاده از داده‌های سری آزمون

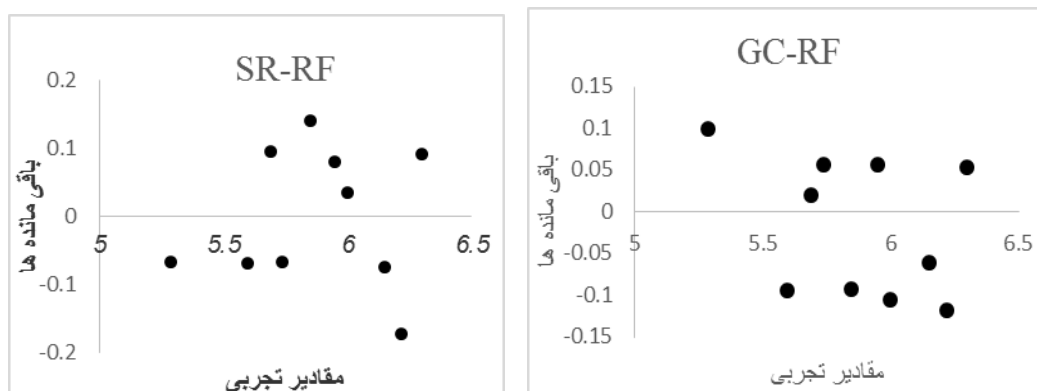
شماره ترکیب	مقدار تجربی (pEC50)	مقدار پیش‌بینی شده		درصد خطا	
		GC-RF	SR - RF	GC- RF	SR - RF
۶	۵/۷۴	۵/۸۰	۵/۸۱	۱/۰۴	۱/۲۱
۱۱	۵/۶۹	۵/۷۱	۵/۵۹	۰/۳۵	-۱/۷۵
۲۳	۵/۶۰	۵/۵۰	۵/۶۷	-۱/۷۸	۱/۲۵
۲۸	۶/۱۵	۶/۰۹	۶/۲۲	-۰/۹۷	۱/۱۳
۳۰	۵/۸۵	۵/۷۶	۵/۷۰	-۱/۵۳	-۲/۳۹
۳۴	۵/۲۹	۵/۳۹	۵/۳۶	۱/۸۹	۱/۳۲
۴۳	۵/۹۵	۶/۰۱	۵/۸۷	۱/۰۰	-۱/۳۴
۴۵	۶/۰۰	۵/۸۹	۵/۹۶	-۱/۸۳	-۰/۶۶
۴۸	۶/۲۲	۶/۱۰	۶/۳۹	-۱/۹۲	۲/۷۳
۵۲	۶/۳۰	۶/۳۵	۶/۲۱	۰/۷۹	-۱/۴۲



شکل ۳-۵: نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی توسط مدل‌های GC-RF و SR-RF برای داده‌های سری آزمون

### ۳-۶-۲- ارزیابی مدل‌های RF با استفاده از نمودار خطای باقی‌مانده

اختلاف مقادیر پیش‌بینی شده و مقادیر تجربی، خطای باقیمانده نامیده می‌شود. توزیع متقارن داده‌ها حول محور افقی (خطای صفر) حاکی از عدم وجود خطای سیستماتیک است. نمودار خطای باقیمانده برحسب مقادیر تجربی، برای مدل‌های ذکر شده در شکل ۳-۶ نشان داده شده است.



شکل ۳-۶: نمودار خطای باقیمانده برحسب مقادیر تجربی برای داده‌های سری آزمون با مدل‌های جنگل تصادفی GC-RF و

SR-RF

### ۳-۶-۳-۳- ارزیابی مدل RF به وسیله روش رد مرحله‌ای تک تک

برای ارزیابی مدل توسط این روش هر بار یک ترکیب از ۵۳ ترکیب به عنوان سری آزمون کنار گذاشته شد و در شرایط بهینه به دست آمده برای سری آموزش با ترکیبات باقی‌مانده مدل‌سازی صورت گرفت. سپس مدل بدست آمده برای پیش‌بینی  $pEC_{50}$  ترکیب کنار گذاشته شده به کار گرفته شد و این فرآیند برای تمام اعضای سری داده‌ها تکرار شد. نتایج حاصل از این روش در جدول ۳-۹ ارائه شده است. نتایج به دست آمده در شکل ۳-۷ نشان می‌دهند که مدل جنگل‌های تصادفی توانایی خوبی برای پیش‌بینی فعالیت ضد ایدز مشتقات پیریدینون دارد. همچنین در شکل ۳-۸ که مقادیر باقی‌مانده‌ها بر حسب مقادیر تجربی  $pEC_{50}$  ترکیبات مورد بحث ترسیم شده است پراکندگی نقاط در دو طرف محور افقی متقارن است که نشان دهنده‌ی عدم وجود خطای سیستماتیک برای بیشتر ترکیبات است.

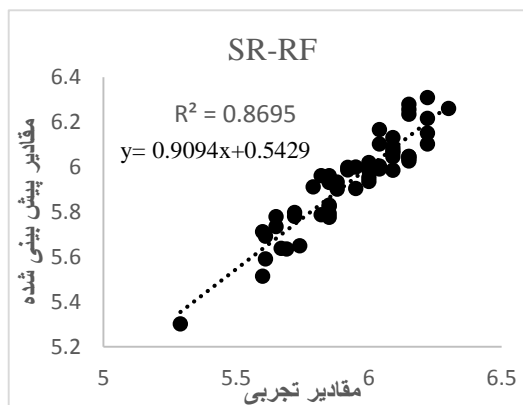
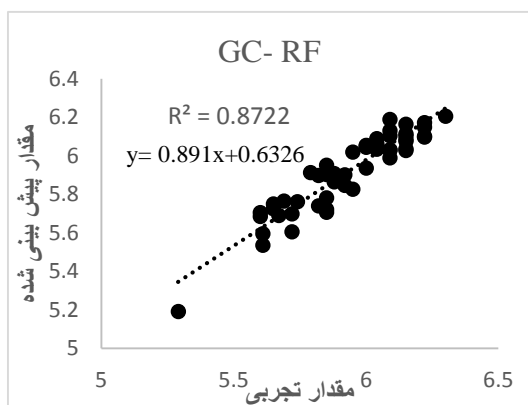
جدول ۳-۹: نتایج حاصل از رد مرحله‌ای تک‌تک برای مدل‌های به دست آمده توسط روش جنگل‌های تصادفی

شماره ترکیب	مقدار تجربی (pEC50)	مقدار پیش‌بینی		درصد خطا	
		GC-RF	SR - RF	GC- RF	SR - RF
۱	۵/۶۷	۵/۶۹	۵/۶۴	۰/۳۵	-۰/۵۲
۲	۶/۰۴	۶/۰۷	۵/۹۹	۰/۴۹	-۰/۸۲
۳	۶/۰۰	۶/۰۵	۵/۹۸	۰/۸۳	-۰/۳۳
۴	۶/۲۲	۶/۱۰	۶/۱۰	-۱/۹۲	-۱/۹۲
۵	۶/۱۵	۶/۱۶	۶/۰۵	۰/۱۶	-۱/۶۲
۶	۵/۷۴	۵/۷۶	۵/۶۵	۰/۳۴	-۱/۵۶
۷	۶/۰۹	۶/۰۳	۶/۰۴	-۰/۹۸	-۰/۸۲
۸	۵/۷۹	۵/۹۱	۵/۹۱	۲/۱۱	۲/۰۷
۹	۶/۰۹	۶/۰۳	۶/۱۳	-۰/۹۸	۰/۶۵
۱۰	۶/۵۲	۶/۰۵	۶/۰۲	۰/۸۳	۰/۳۳
۱۱	۵/۶۹	۵/۷۷	۵/۶۳	۱/۴۰	-۱/۰۵
۱۲	۵/۹۵	۵/۸۳	۵/۹۰	-۲/۰۱	-۰/۸۴
۱۳	۵/۹۲	۵/۸۵	۶/۰۰	-۱/۱۸	۱/۳۵
۱۴	۶/۱۵	۶/۰۳	۶/۰۴	-۱/۹۵	-۱/۷۸
۱۵	۵/۶۱	۵/۶۰	۵/۶۹	-۰/۱۷	۱/۴۲
۱۶	۵/۶۰	۵/۷۱	۵/۷۱	۱/۹۶	۱/۹۶
۱۷	۶/۱۵	۶/۱۱	۶/۰۳	-۰/۶۵	-۱/۹۵
۱۸	۶/۲۲	۶/۱۷	۶/۲۲	-۰/۸۰	۰/۰۰
۱۹	۵/۸۵	۵/۷۲	۵/۷۹	-۲/۲۲	-۱/۰۲
۲۰	۵/۸۵	۵/۹۵	۵/۹۶	۱/۷۰	-۱/۳۶
۲۱	۶/۰۹	۶/۱۰	۶/۰۶	۰/۱۶	-۰/۴۹
۲۲	۵/۶۱	۵/۵۴	۵/۵۹	-۱/۲۴	-۰/۳۵
۲۳	۵/۶۰	۵/۶۹	۵/۵۱	۱/۶۰	-۱/۶۰
۲۴	۵/۸۵	۵/۹۰	۵/۸۳	۰/۸۵	-۰/۳۴
۲۵	۵/۸۵	۵/۹۲	۵/۹۶	۱/۱۹	۱/۸۸
۲۶	۶/۲۲	۶/۱۰	۶/۳۱	-۱/۹۲	۱/۴۴
۲۷	۵/۶۵	۵/۷۲	۵/۷۸	۱/۲۳	۲/۳۰
۲۸	۶/۱۵	۶/۱۱	۶/۲۸	-۰/۶۵	۲/۱۱
۲۹	۶/۰۴	۶/۰۳	۶/۱۰	-۰/۱۶	۰/۹۹
۳۰	۵/۸۵	۵/۷۸	۵/۸۲	-۱/۱۹	-۰/۵۱
۳۱	۵/۸۵	۵/۷۱	۵/۹۳	-۲/۳۹	۱/۳۶

ادامه جدول (۹-۳)

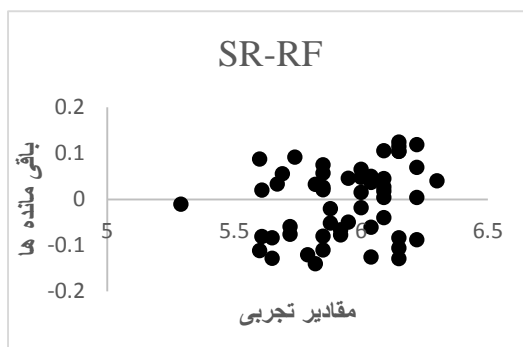
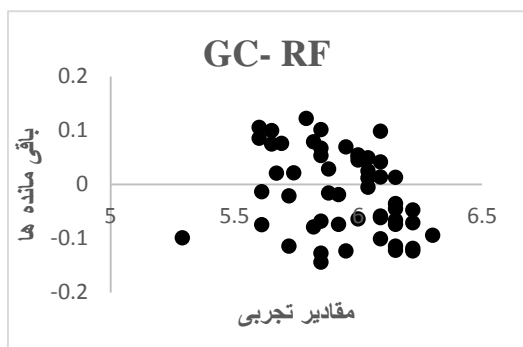
شماره ترکیب	مقدار تجربی (pEC50)	مقدار پیش بینی		درصد خطا	
		GC-RF	SR - RF	GC- RF	SR - RF
۳۲	۵/۹۲	۵/۹۰	۵/۹۹	-۰/۳۳	۱/۱۸
۳۳	۶/۱۵	۶/۰۸	۶/۲۳	-۱/۱۳	۱/۳۰
۳۴	۵/۲۹	۵/۱۹	۵/۳۰	-۱/۸۹	۰/۱۸
۳۵	۵/۸۲	۵/۹۰	۵/۹۶	۱/۳۷	۲/۴۰
۳۶	۶/۰۹	۶/۱۳	۶/۰۹	۰/۶۵	۰/۰۰
۳۷	۵/۶۵	۵/۷۵	۵/۷۳	۱/۷۶	۱/۴۱
۳۸	۵/۷۲	۵/۶۱	۵/۸۰	-۱/۹۲	۱/۳۹
۳۹	۵/۷۲	۵/۷۰	۵/۷۸	-۰/۳۴	۱/۰۴
۴۰	۶/۰۴	۶/۰۵	۶/۱۷	۰/۱۶	۲/۱۵
۴۱	۶/۰۴	۶/۰۹	۶/۰۰	۰/۸۲	-۰/۶۶
۴۲	۶/۱۵	۶/۰۴	۶/۲۶	-۱/۷۸	۱/۷۸
۴۳	۵/۹۵	۶/۰۱	۶/۰۰	۱/۱۷	۰/۸۴
۴۴	۶/۰۹	۵/۹۹	۵/۹۸	-۱/۶۴	-۱/۸۰
۴۵	۶	۵/۹۴	۵/۹۵	-۱/۰۰	-۰/۸۳
۴۶	۶/۰۹	۶/۱۹	۶/۰۷	۱/۶۴	-۰/۳۲
۴۷	۵/۸۲	۵/۷۴	۵/۷۹	-۱/۳۷	-۰/۵۱
۴۸	۶/۲۲	۶/۱۵	۶/۱۵	-۱/۱۲	-۱/۱۲
۴۹	۶/۱۵	۶/۰۸	۶/۰۵	-۱/۱۳	-۱/۶۲
۵۰	۵/۸۸	۵/۸۶	۵/۹۰	-۰/۳۴	۰/۳۴
۵۱	۶	۶/۰۵	۵/۹۳	۰/۸۳	-۱/۱۶
۵۲	۶/۳۰	۶/۲۱	۶/۲۶	-۱/۴۲	-۰/۶۳
۵۳	۵/۸۸	۵/۹۱	۵/۹۳	۰/۵۱	۰/۸۵





شکل ۳-۷: نمودار مقادیر پیش‌بینی شده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک

توسط مدل‌های جنگل‌های تصادفی (GC-RF و SR-RF) بر حسب مقادیر تجربی

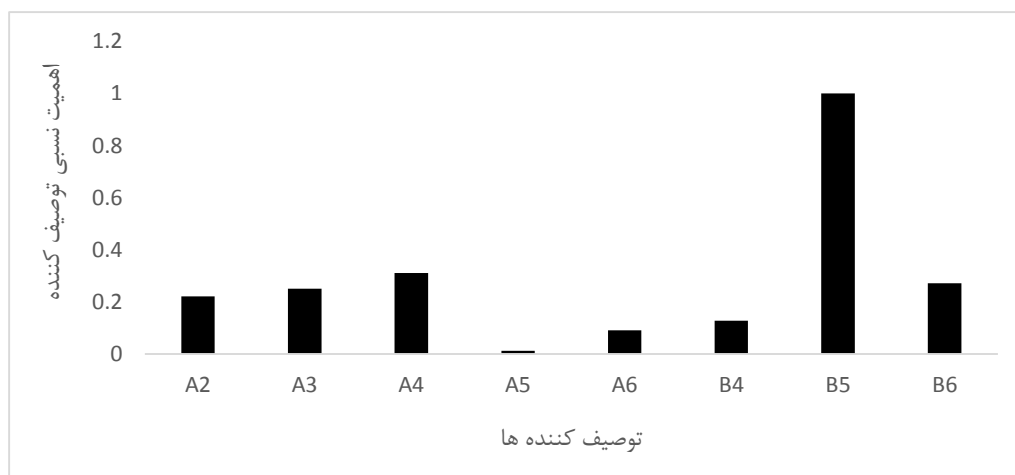


شکل ۳-۸: نمودار مقادیر خطای باقی‌مانده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک توسط مدل‌های

جنگل‌های تصادفی (GC-RF و SR-RF) بر حسب مقادیر تجربی

### ۳-۶-۴- انتخاب توصیف کننده‌های مهم با روش جنگل‌های تصادفی

روش جنگل‌های تصادفی علاوه بر اینکه به عنوان یک روش مدل‌سازی به کار می‌رود به عنوان یک روش انتخاب متغیر نیز به کار برده می‌شود. پس از ساخت مدل جنگل‌های تصادفی و بهینه کردن پارامترهای موثر، الگوریتم جنگل‌های تصادفی با استفاده از مقادیر بهینه صد بار تکرار شد و در هر بار میزان کاهش MSE به عنوان مقدار اهمیت توصیف کننده‌ها محاسبه گردید. سرانجام میانگین مقادیر اهمیت هر توصیف کننده به عنوان شاخص نهایی در نظر گرفته شد و سپس اهمیت نسبی هر توصیف کننده نسبت به توصیف کننده‌ای که بیشترین اهمیت را داراست، تعیین شد. شکل ۳-۹ میزان اهمیت نسبی توصیف کننده‌ها را نمایش می‌دهد و همانطور که مشاهده می‌شود بعضی از توصیف کننده‌ها دارای اهمیت نسبی کمی بوده و تنها ۴ توصیف کننده اهمیت نسبی قابل توجهی دارند که در جدول ۳-۱۰ این توصیف کننده‌ها همراه با مقدار اهمیت نسبی و ضریب اثر آن‌ها گردآوری شده‌اند.



شکل ۳-۹: نمودار اهمیت نسبی توصیف کننده‌های سهم گروه

جدول ۳-۱۰: توصیف کننده‌های برتر با استفاده از روش جنگل‌های تصادفی

No	علامت	مفهوم توصیف کننده	اهمیت نسبی
۱	A <sub>4</sub>	حضور یا عدم حضور استخلاف در موقعیت پارا حلقه A	۰/۳۱۱۳
۲	A <sub>6</sub>	حضور یا عدم حضور استخلاف در موقعیت ارتو حلقه A	۰/۰۹۰۶
۳	B <sub>5</sub>	حضور یا عدم حضور استخلاف در موقعیت متا حلقه B	۱/۰۰۰۰
۴	B <sub>6</sub>	حضور یا عدم حضور استخلاف در موقعیت ارتو حلقه B	۰/۲۷۱۱

### ۳-۷- مدل‌سازی شبکه عصبی مصنوعی

#### ۳-۷-۱- مدل‌سازی شبکه عصبی مصنوعی با توصیف‌کننده‌های سهم

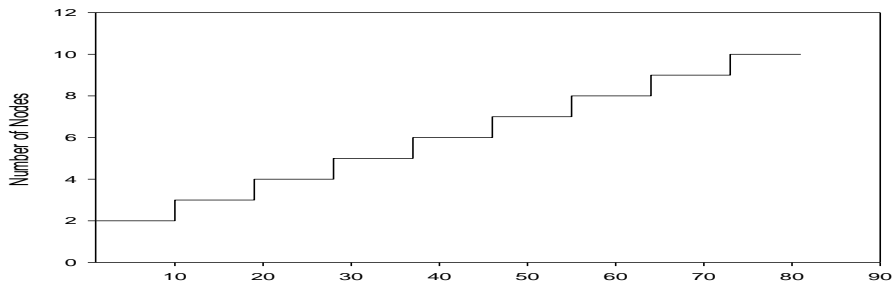
##### گروه (GC-ANN)

یکی از راه‌های یافتن رابطه غیرخطی بین متغیرهای مستقل و متغیر وابسته، استفاده از شبکه عصبی مصنوعی برای مدل‌سازی می‌باشد. شبکه عصبی انتخاب شده در این مطالعه یک شبکه پیشرو با الگوریتم آموزشی پس انتشار می‌باشد که الگوریتم آن با برنامه متلب نوشته شده است در فرآیند بهینه‌سازی پارامترهای شبکه، سری داده‌ها به سه بخش سری آزمون (۸ ترکیب) و سری ارزیابی (۸ ترکیب) و همچنین سری آموزش (۳۷ ترکیب) تقسیم شدند. ۸ توصیف‌کننده‌ای که از روش سهم گروه به دست آمدند به عنوان ورودی به شبکه عصبی مصنوعی داده شدند. خروجی شبکه نیز فعالیت متناظر با ورودی‌ها را نشان می‌دهد. برای آموزش شبکه، دو الگوریتم آموزشی لونبرگ-مارکوات (Im) و تنظیم بایزین (br)، و دوتابع انتقال تانژانت سیگموئید (tan sig) و لگاریتم سیگموئید (log sig) مورد استفاده قرار گرفتند که از ترکیب این دو الگوریتم و دو تابع انتقال، چهار نوع شبکه عصبی مختلف طراحی گردید. در هر یک از این شبکه‌ها باید تعداد ورودی‌ها و تعداد نرون‌های لایه پنهان (گره) و تعداد دور آموزشی بهینه گردند که مقدار بهینه هر یک از این پارامترها زمانی است که میانگین مربع خطا (MSE) برای سری ارزیابی دارای کمترین مقدار باشد. برای بهینه‌سازی تعداد ورودی‌های

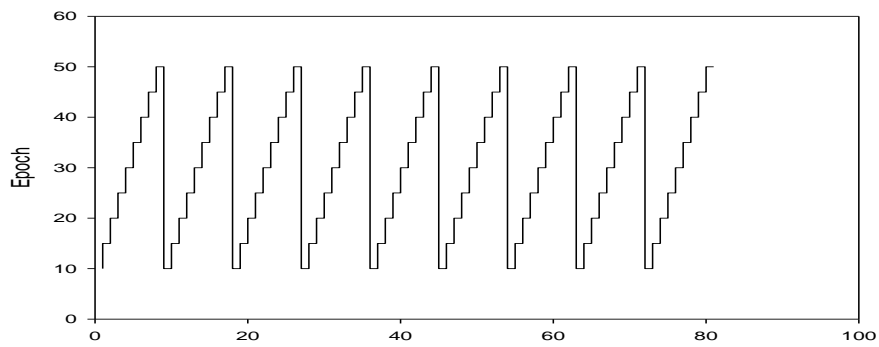
شبکه، تعداد نرون‌های لایه‌ی پنهان و تعداد دور آموزشی از روش بهینه‌سازی همزمان استفاده گردید. برای این منظور هر شبکه با تعداد ورودی ۸ و تعداد گره از ۲ تا ۱۰ و تعداد دور آموزشی از ۱۰ تا ۵۰ با گام ۵ به طور همزمان آموزش داده شد. در روند بهینه‌سازی فوق، به حداقل رساندن میانگین خطاهای سری ارزیابی به عنوان معیار انتخاب شد، این کار موجب می‌شود بتوانیم نقش تمام پارامترها را به طور همزمان و یکجا بررسی نماییم. بخشی از روند تغییرات پارامترهای شبکه در حین بهینه‌سازی همزمان پارامترها به همراه مقادیر MSE بدست آمده، به صورت نموداری برحسب یک بردار مرجع فرضی در شکل ۳-۱۰ تا شکل ۳-۱۳ آمده است.

با توجه به نتایج حاصله برای توصیف‌کننده‌ها و تعداد مختلف گره‌ها و دوره‌های آموزش و همچنین توابع متفاوت آموزش و انتقال، بهترین شبکه‌های بدست آمده براساس کمترین مقدار MSE در جدول ۳-۱۱ خلاصه شده است.

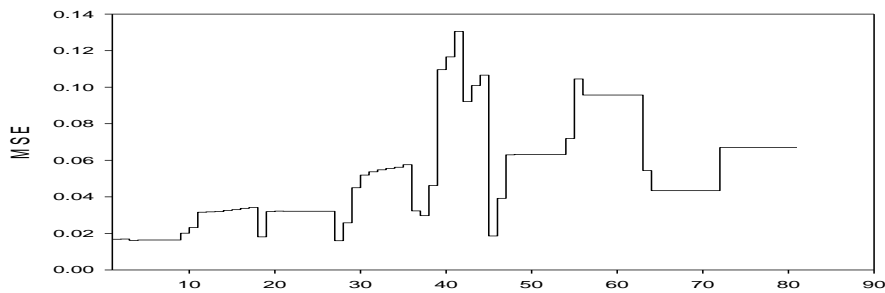
با توجه به جدول ۳-۱۱ شبکه عصبی با ۸ توصیف‌کننده در لایه ورودی، سه نرون در لایه مخفی و تعداد دور آموزشی ۴۵ با تابع آموزش لونیبرگ\_مارکوات و تابع انتقال تانژانت سیگموئیدی کمترین MSE را نسبت به سایر شبکه‌ها نشان می‌دهد. بنابراین شبکه‌ای با این ساختار به عنوان بهترین مدل ANN برای مدل‌سازی داده‌ها در نظر گرفته شد.



(الف)



(ب)

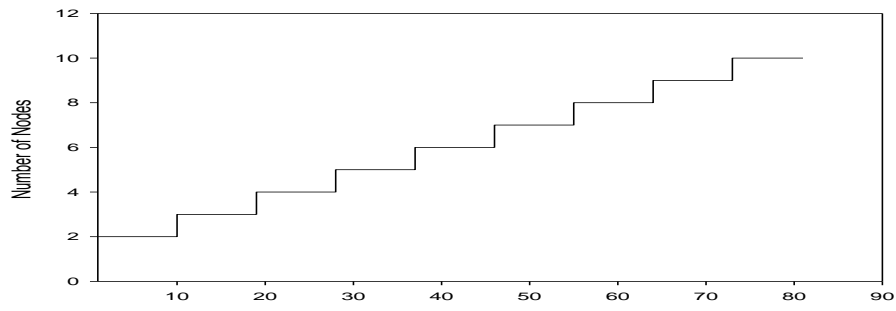


بردار مرجع

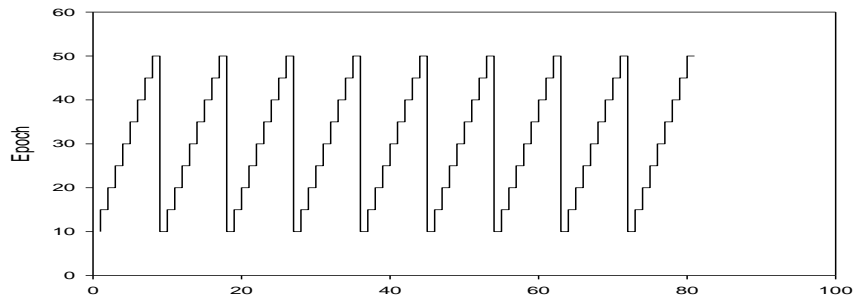
(ج)

شکل ۳-۱۰: نمودارهای الف) تعداد گره‌ها، ب) تعداد دوره‌های آموزش و ج) مقادیر MSE برای تابع آموزش لونیبرگ-

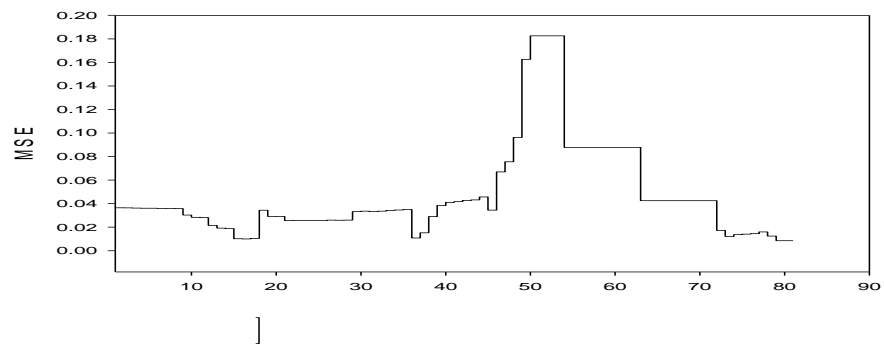
مارکوات به همراه تابع انتقال لگاریتم سیگموئید



(الف)



(ب)

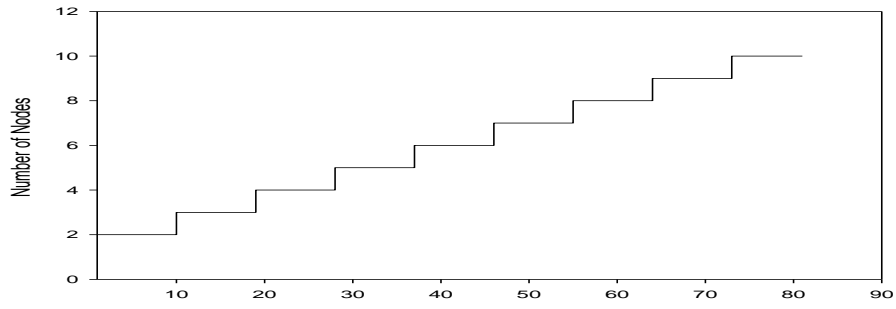


بردار مرجع

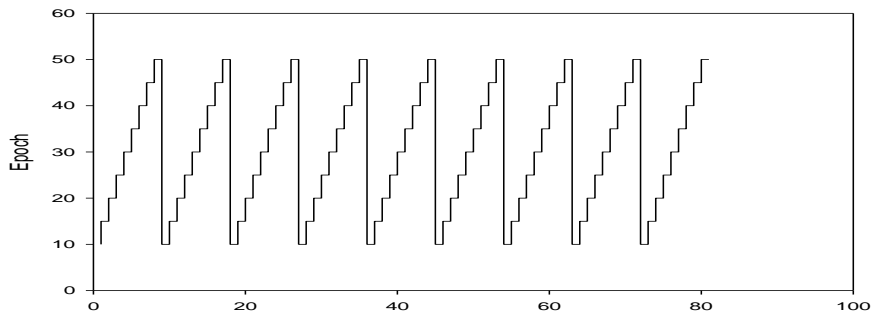
(ج)

شکل ۳-۱۱: نمودارهای الف) تعداد گره‌ها، ب) تعداد دوره‌های آموزش و ج) مقادیر MSE برای تابع آموزش لونیبرگ-

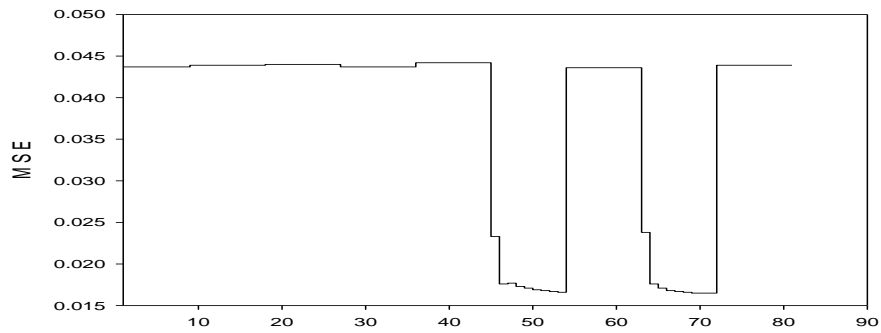
مارکوات به همراه تابع انتقال تانژانت سیگموئید



(الف)



(ب)

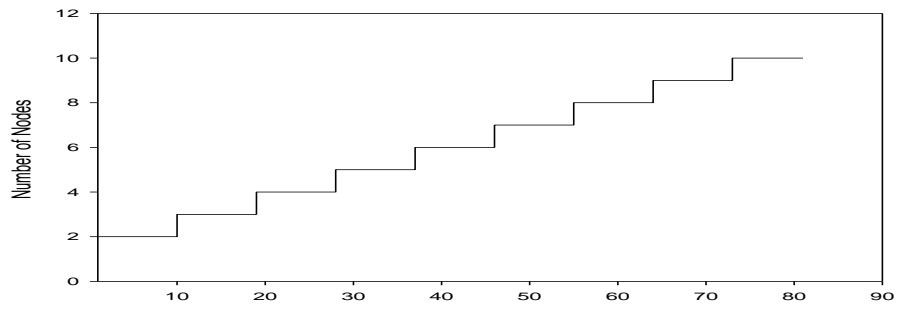


بردار مرجع

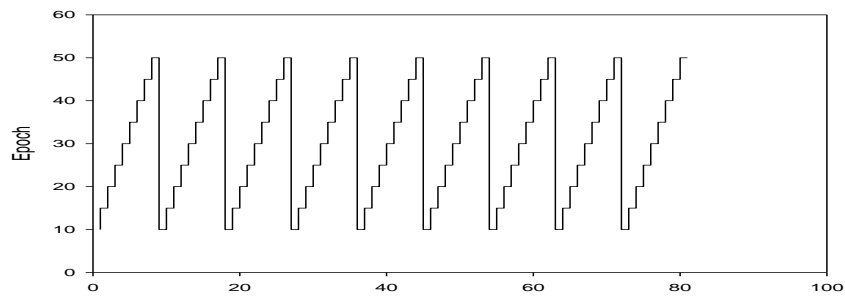
(ج)

شکل ۳-۱۲: نمودارهای الف) تعداد گره‌ها، ب) تعداد دوره‌های آموزش و ج) مقادیر MSE برای تابع آموزش تنظیم

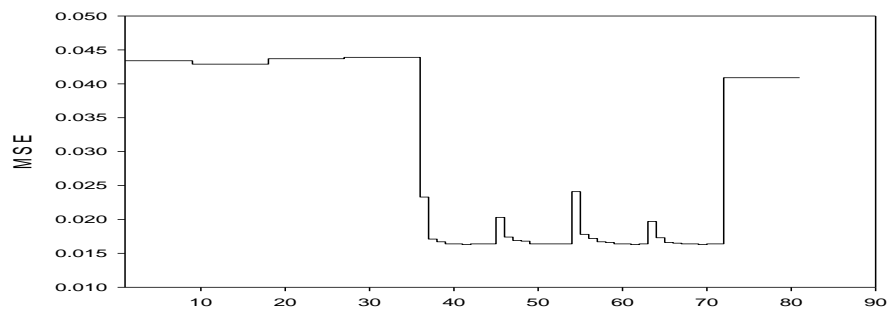
بایزین به همراه تابع انتقال لگاریتم سیگموئید



(الف)



(ب)



بردار مرجع

(ج)

شکل ۳-۱۳: نمودارهای الف) تعداد گره‌ها، ب) تعداد دوره‌های آموزش و ج) مقادیر  $MSE$  برای تابع آموزش تنظیم

بایزین به همراه تابع انتقال تانژانت سیگموئید



جدول ۳-۱۱: توابع و پارامترهای شبکه‌های بهینه GC-ANN

MSE	تعداد دور آموزش	تعداد گره	تابع انتقال	تابع آموزش	تعداد توصیف کننده
۰/۰۱۶۰	۱۰	۵	لگاریتم-سیگموئید	لونبرگ-مارکوات	۸
۰/۰۱۰۰	۴۵	۳	تانژانت-سیگموئید	لونبرگ-مارکوات	۸
۰/۰۱۶۵	۴۰	۹	لگاریتم-سیگموئید	تنظیم بایزین	۸
۰/۰۱۶۳	۳۵	۶	تانژانت-سیگموئید	تنظیم بایزین	۸

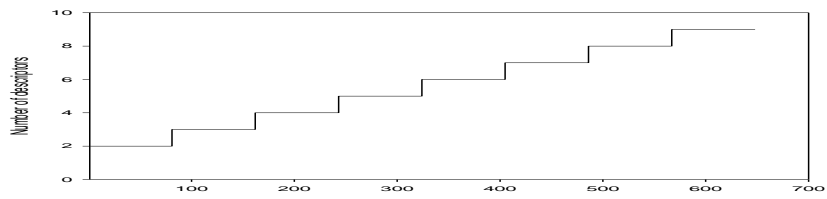
### ۳-۷-۲- مدل‌سازی شبکه عصبی مصنوعی با توصیف کننده‌های

#### انتخاب شده به وسیله روش رگرسیون مرحله‌ای (SR - ANN)

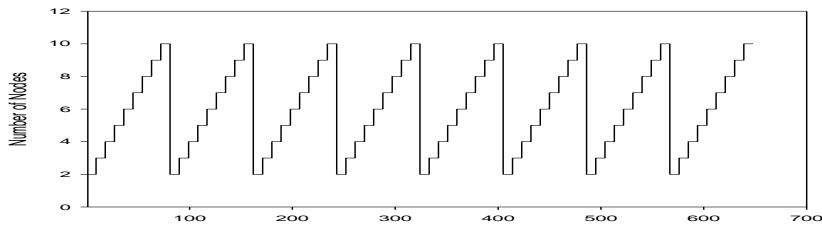
ابتدا مقادیر ۹ توصیف کننده انتخاب شده با روش SR (رگرسیون مرحله‌ای) به عنوان متغیرهای ورودی و فعالیت متناظر آن‌ها به عنوان متغیر هدف در نظر گرفته شد تا پاسخ شبکه با آن‌ها سنجیده شود. برای بدست آوردن بهترین مدل و کمترین خطا، پارامترهای موثر (تابع انتقال، تابع آموزش، تعداد متغیرهای ورودی شبکه، تعداد گره‌ها در لایه مخفی و تعداد دورهای آموزشی) به‌طور همزمان بهینه شدند. تابع کارایی شبکه نیز میانگین مربع خطا (MSE) بود.

در فرآیند بهینه‌سازی پارامترهای شبکه، سری داده‌ها به سه بخش سری آزمون (۸ ترکیب) و سری ارزیابی (۸ ترکیب) و همچنین سری آموزش (۳۷ ترکیب) تقسیم شدند. سری ارزیابی برای محاسبه گرادیان و به روز کردن اوزان و بایاس‌های شبکه استفاده شد. از آنجایی که در بیشتر موارد به نظر می‌رسد که یک لایه پنهان در ساختار شبکه عصبی مناسب باشد، در این تحقیق نیز از یک لایه پنهان استفاده شده است [۴۰]. برای بهینه‌سازی تعداد متغیرهای ورودی شبکه، تعداد نرون‌های لایه پنهان و تعداد دور آموزشی از روش بهینه‌سازی همزمان استفاده گردید. برای یافتن مقدار بهینه پارامترهای موثر، مقادیر مختلفی از آن‌ها در ساختار شبکه قرار داده شد. برای این منظور هر شبکه با تعداد ورودی ۲ تا ۹ و با دو الگوریتم آموزشی لونبرگ-مارکوات (train lm) و تنظیم بایزین (train br) و دو تابع انتقال لگاریتم سیگموئید و تانژانت سیگموئید و تعداد گره از ۲ تا ۱۰ و تعداد دورهای

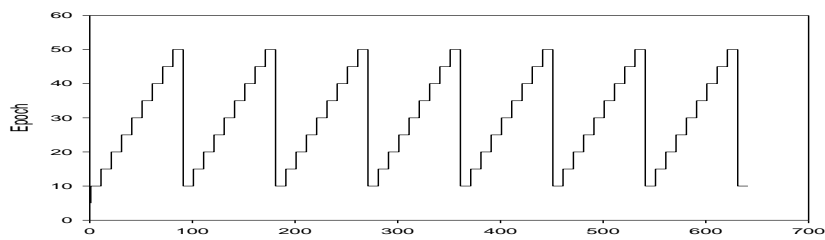
آموزش از ۱۰ تا ۵۰ با گام ۵، به طور همزمان آموزش داده شد. در روند بهینه‌سازی فوق، به حداقل رساندن میانگین خطاهای سری ارزیابی به عنوان معیار انتخاب شد. بخشی از روند تغییرات پارامترهای شبکه در حین بهینه‌سازی همزمان پارامترها به همراه مقادیر  $MSE$  بدست آمده، به صورت نموداری برحسب یک بردار مرجع فرضی در شکل ۳-۱۴ تا شکل ۳-۱۷ آمده است. نتیجه بهینه‌سازی این پارامترها در جدول ۳-۱۲ گردآوری شده است.



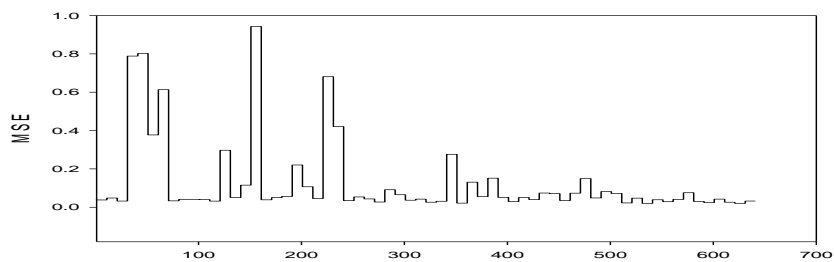
(الف)



(ب)



(ج)

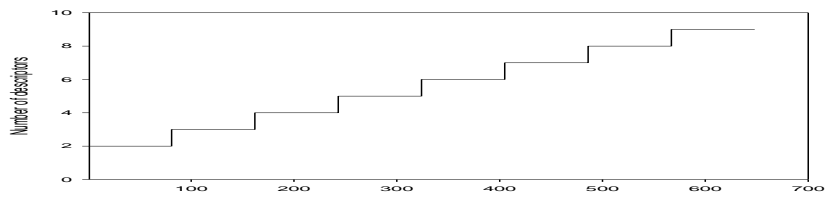


بردار مرجع

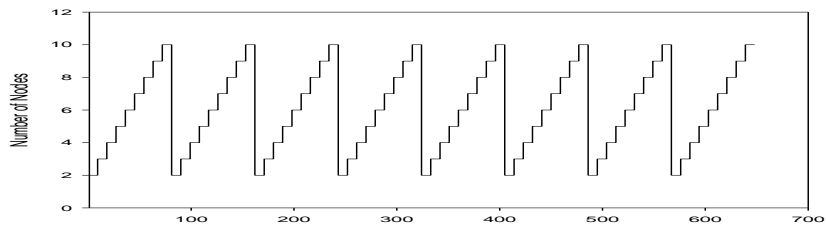
(د)

شکل ۳-۱۴: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دورهای آموزش و د) مقادیر MSE برای

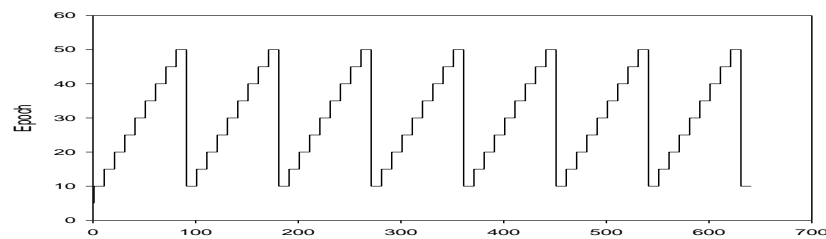
تابع آموزش لونیبرگ-مارکوات به همراه تابع انتقال لگاریتم سیگموئید



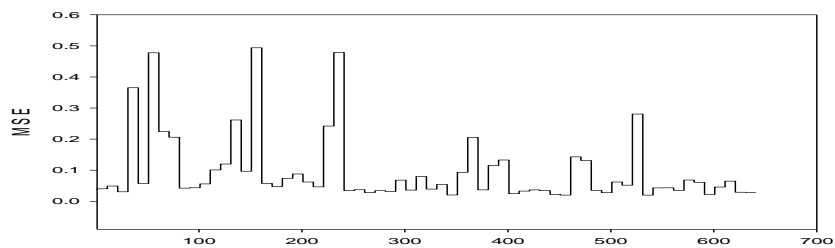
(الف)



(ب)



(ج)

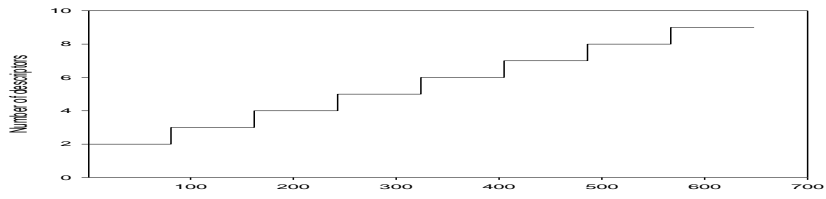


بردار مرجع

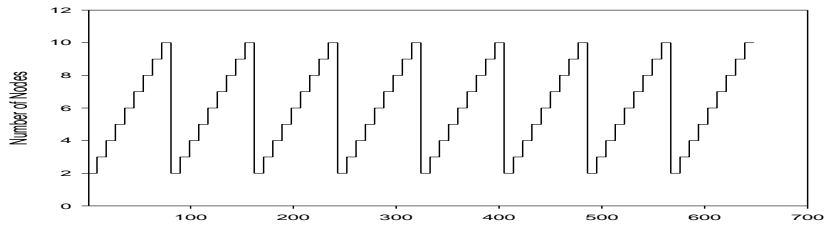
(د)

شکل ۳-۱۵: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای

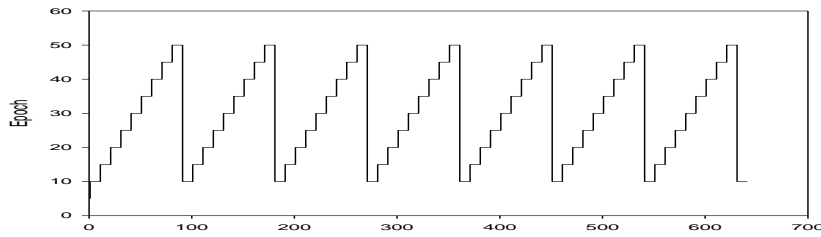
تابع آموزش لونیبرگ-مارکوات به همراه تابع انتقال تانژانت سیگموئید



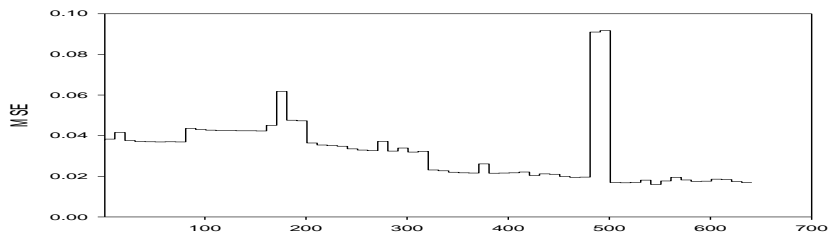
(الف)



(ب)



(ج)

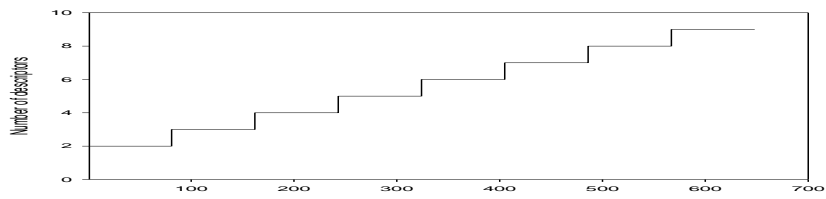


بردار مرجع

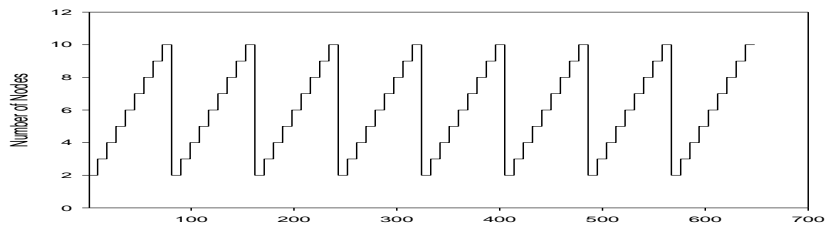
(د)

شکل ۳-۱۶: نمودارهای الف) تعداد توصیف کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای

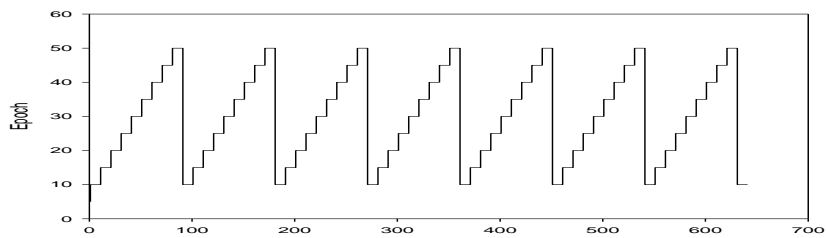
تابع آموزش تنظیم بایزین به همراه تابع انتقال لگاریتم سیگموئید



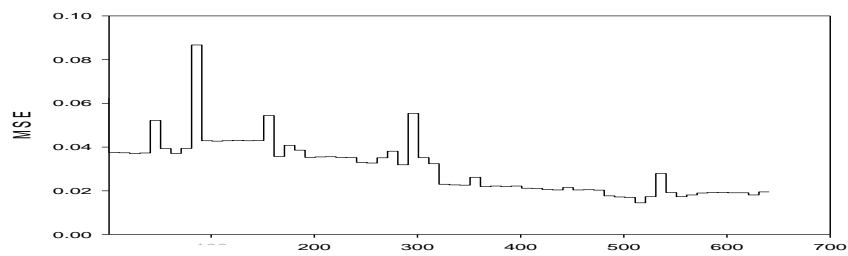
(الف)



(ب)



(ج)



بردار مرجع

(د)

شکل ۳-۱۷: نمودارهای الف) تعداد توصیف کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای

تابع آموزش تنظیم بایزین به همراه تابع انتقال تانژانت سیگموئید

جدول ۳-۱۲: توابع و پارامترهای شبکه‌های بهینه SR-ANN

MSE	تعداد دور آموزش	تعداد گره	تابع انتقال	تابع آموزش	تعداد توصیف کننده
۰/۰۱۶۱	۱۰	۶	لگاریتم-سیگموئید	لونبرگ-مارکوات	۹
۰/۰۰۹۶	۱۵	۷	تانژانت-سیگموئید	لونبرگ-مارکوات	۸
۰/۰۱۵۶	۱۵	۶	لگاریتم-سیگموئید	تنظیم بایزین	۸
۰/۰۱۴۵	۴۵	۵	تانژانت-سیگموئید	تنظیم بایزین	۸

با توجه به نتایج بدست آمده، الگوریتم آموزشی لونبرگ-مارکوات و تابع انتقال تانژانت سیگموئید، ۸ توصیف کننده ورودی و ۷ نرون در لایه مخفی با تعداد دور آموزشی ۱۵ کمترین MSE را نشان می‌دهد. بنابراین این شبکه برای مدل‌سازی در نظر گرفته شد.

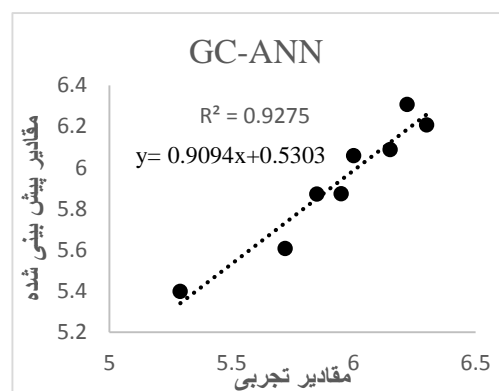
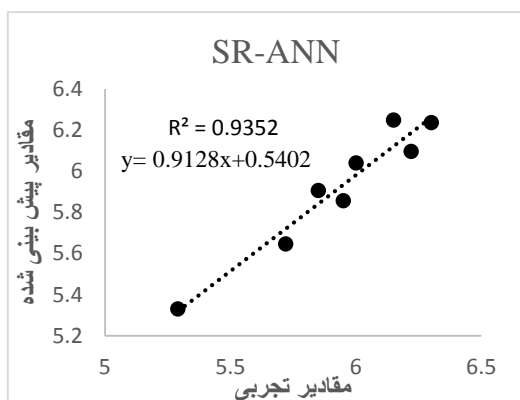
### ۳-۷-۳- ارزیابی مدل‌های شبکه عصبی مصنوعی

#### ۳-۷-۳-۱- ارزیابی مدل‌های ANN با استفاده از داده‌های سری آزمون

جدول ۱۳-۳ نتایج پیش‌بینی حاصل از ارزیابی مدل‌های شبکه عصبی مصنوعی با استفاده از داده‌های سری آزمون و شکل ۱۸-۳ نیز نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی را برای داده‌های سری آزمون نشان می‌دهد.

جدول ۱۳-۳: نتایج حاصل از ارزیابی مدل‌های شبکه عصبی مصنوعی با استفاده از داده‌های سری آزمون

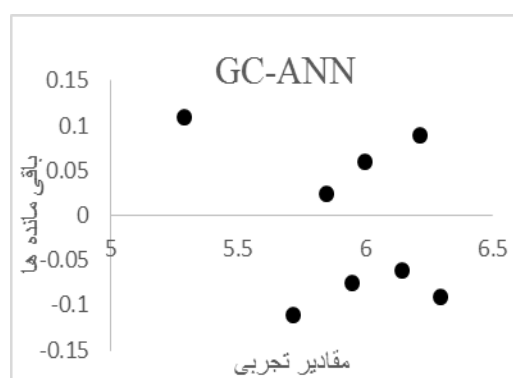
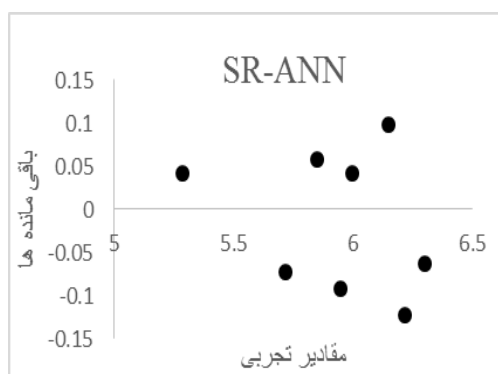
شماره ترکیب	مقدار تجربی	مقدار پیش‌بینی		درصد خطا	
		GC-ANN	SR-ANN	GC-ANN	SR-ANN
۱۰	۶/۰۰	۶/۰۶	۵/۸۷	۱/۰۰	-۲/۱۶
۱۲	۵/۹۵	۵/۸۷	۶/۰۲	-۱/۳۴	۱/۱۷
۱۸	۶/۲۲	۶/۳۱	۶/۱۹	۱/۴۴	-۰/۴۸
۲۰	۵/۸۵	۵/۸۷	۵/۹۳	۰/۳۴	۱/۳۶
۳۳	۶/۱۵	۶/۰۹	۶/۱۱	-۰/۹۷	-۰/۶۵
۳۴	۵/۲۹	۵/۴۰	۵/۳۵	۲/۰۷	۱/۱۳
۳۹	۵/۷۲	۵/۶۱	۵/۸۲	-۱/۹۲	۱/۷۴
۵۲	۶/۳۰	۶/۲۱	۶/۳۹	-۱/۴۲	۱/۴۲



شکل ۳-۱۸: نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی را برای داده‌های سری آزمون

### ۳-۷-۳-۲- ارزیابی مدل‌های ANN با استفاده از نمودار خطای باقی‌مانده

شکل ۳-۱۹ نمودار خطای باقیمانده برحسب مقادیر تجربی، برای مدل‌های ذکر شده را نشان می‌دهد. توزیع متقارن داده‌ها حول محور افقی (خطای صفر) حاکی از عدم وجود خطای سیستماتیک است.



شکل ۳-۱۹: نمودار خطای باقیمانده برحسب مقادیر تجربی برای داده‌های سری آزمون



### ۳-۳-۷-۳- ارزیابی مدل‌های ANN توسط روش رد مرحله‌ای تک تک

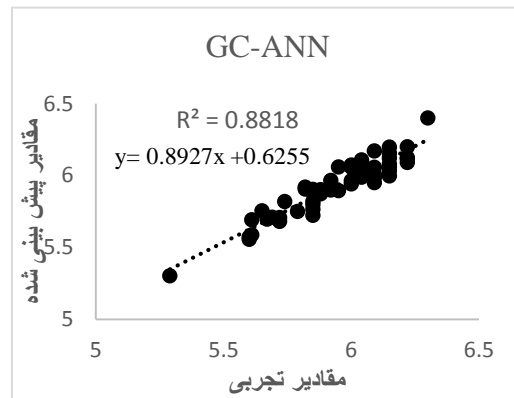
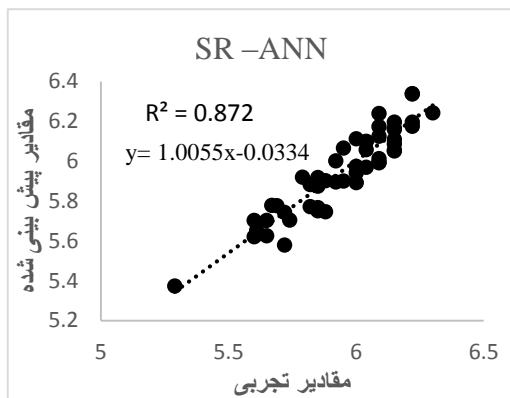
برای ارزیابی مدل‌های پیشنهادی از روش رد مرحله‌ای تک تک نیز استفاده گردید. نتایج حاصل از این روش در جدول ۳-۱۴ ارائه شده است. نتایج به دست آمده در شکل ۳-۲۰ نشان می‌دهند که مدل شبکه عصبی مصنوعی توانایی خوبی برای پیش‌بینی فعالیت ضد ایدز مشتقات پیریدینون دارد. همچنین شکل ۳-۲۱ که مقادیر باقی‌مانده‌ها بر حسب مقادیر تجربی pEC<sub>50</sub> ترکیبات مورد بحث ترسیم شده است دارای تقارن پراکندگی نقاط در دو طرف محور افقی است که نشان دهنده‌ی عدم وجود خطای سیستماتیک برای بیشتر ترکیبات است.

جدول ۳-۱۴: نتایج حاصل از ارزیابی مدل‌های (GC-ANN و SR-ANN) به دست آمده توسط روش رد مرحله‌ای تک تک

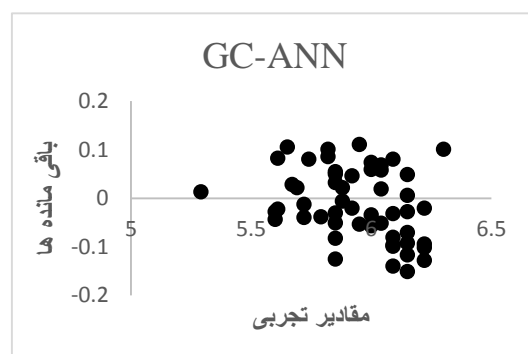
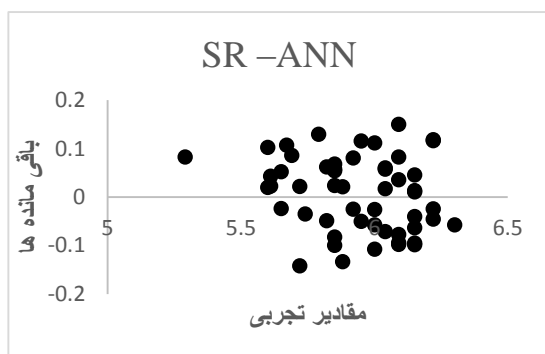
شماره ترکیب	مقدار تجربی (pEC <sub>50</sub> )	مقدار پیش‌بینی		درصد خطا	
		GC-ANN	SR -ANN	GC-ANN	SR -ANN
۱	۵/۶۷	۵/۷۰	۵/۵۴	۰/۵۲	-۲/۲۹
۲	۶/۰۴	۶/۱۱	۵/۹۴	۱/۱۵	-۱/۶۵
۳	۶/۰۰	۶/۰۷	۵/۹۴	۱/۱۶	-۱/۰۰
۴	۶/۲۲	۶/۱۳	۶/۲۲	-۱/۴۴	۰/۰۰
۵	۶/۱۵	۶/۰۰	۶/۲۴	-۲/۴۳	-۱/۴۶
۶	۵/۷۴	۵/۸۲	۵/۶۷	۱/۳۹	-۱/۲۱
۷	۶/۰۹	۶/۰۱	۵/۹۹	-۱/۳۱	-۱/۶۴
۸	۵/۷۹	۵/۷۵	۵/۷۰	-۰/۶۹	۱/۵۵
۹	۶/۰۹	۵/۹۹	۶/۰۶	-۱/۶۴	-۰/۴۹
۱۰	۶/۰۰	۶/۰۶	۵/۹۵	۱/۰۰	-۰/۸۳
۱۱	۵/۶۹	۵/۷۱	۵/۷۲	۰/۳۵	۰/۵۲
۱۲	۵/۹۵	۵/۹۰	۵/۸۷	-۰/۸۴	-۱/۳۴
۱۳	۵/۹۲	۵/۹۰	۶/۰۱	-۰/۳۳	۱/۵۲
۱۴	۶/۱۵	۶/۱۶	۶/۰۹	۰/۱۶	-۰/۹۷
۱۵	۵/۶۱	۵/۶۹	۵/۵۳	۱/۴۲	-۱/۴۲
۱۶	۵/۶۰	۵/۵۷	۵/۵۷	-۰/۵۳	-۰/۵۳
۱۷	۶/۱۵	۶/۲۰	۶/۲۰	۰/۸۱	-۱/۸۱
۱۸	۶/۲۲	۶/۰۹	۶/۲۷	-۲/۰۹	۰/۸۰
۱۹	۵/۸۵	۵/۷۷	۵/۹۷	-۱/۳۶	۲/۰۵

ادامه جدول (۳-۱۴)

شماره ترکیب	مقدار تجربی (pEC50)	مقدار پیش‌بینی		درصد خطا	
		GC-ANN	SR-ANN	GC-ANN	SR-ANN
۲۰	۵/۸۵	۵/۸۲	۵/۹۴	-۰/۵۱	۱/۵۳
۲۱	۶/۰۹	۵/۹۹	۶/۰۰	-۱/۶۴	-۱/۴۷
۲۲	۵/۶۱	۵/۵۹	۵/۶۹	-۰/۳۵	۱/۴۲
۲۳	۵/۶۰	۵/۵۶	۵/۶۸	-۰/۷۱	۱/۴۲
۲۴	۵/۸۵	۵/۸۰	۵/۷۳	-۰/۸۵	-۲/۰۵
۲۵	۵/۸۵	۵/۹۰	۵/۷۲	۰/۸۵	-۲/۲۲
۲۶	۶/۲۲	۶/۱۲	۶/۳۳	-۱/۶۰	۱/۷۸
۲۷	۵/۶۵	۵/۹۰	۵/۵۳	۰/۸۵	-۲/۱۲
۲۸	۶/۱۵	۶/۰۳	۶/۱۷	-۱/۹۵	۰/۳۲
۲۹	۶/۰۴	۶/۱۰	۶/۱۳	۰/۹۹	۱/۴۹
۳۰	۵/۸۵	۵/۷۲	۵/۸۰	-۲/۲۲	-۰/۸۵
۳۱	۵/۸۵	۵/۸۸	۵/۷۱	۰/۵۱	-۲/۳۹
۳۲	۵/۹۲	۵/۹۷	۶/۰۳	۰/۸۴	۱/۸۵
۳۳	۶/۱۵	۶/۰۸	۶/۱۰	-۱/۱۳	-۰/۸۱
۳۴	۵/۲۹	۵/۳۰	۵/۲۷	۰/۱۸	-۰/۳۷
۳۵	۵/۸۲	۵/۹۱	۵/۹۰	۱/۵۴	۱/۳۷
۳۶	۶/۰۹	۶/۰۶	۶/۱۵	-۰/۴۹	۰/۹۸
۳۷	۵/۶۵	۵/۷۶	۵/۷۳	۱/۹۴	۱/۴۱
۳۸	۵/۷۲	۵/۶۸	۵/۸۰	-۰/۶۹	۱/۳۹
۳۹	۵/۷۲	۵/۷۱	۵/۷۹	-۰/۱۷	۱/۲۲
۴۰	۶/۰۴	۶/۰۶	۶/۱۳	۰/۳۳	۱/۴۹
۴۱	۶/۰۴	۵/۹۹	۶/۰۹	-۰/۸۲	۰/۸۲
۴۲	۶/۱۵	۶/۰۶	۶/۰۵	-۱/۴۶	-۱/۶۲
۴۳	۵/۹۵	۶/۰۶	۶/۰۱	۱/۸۴	۱/۰۰
۴۴	۶/۰۹	۵/۹۵	۶/۰۷	-۲/۲۹	-۰/۳۲
۴۵	۶	۵/۹۷	۶/۰۷	-۰/۵۰	-۱/۱۶
۴۶	۶/۰۹	۶/۱۷	۶/۱۵	۱/۳۱	۰/۹۸
۴۷	۵/۸۲	۵/۹۲	۵/۷۷	۱/۷۱	-۰/۸۵
۴۸	۶/۲۲	۶/۲۰	۶/۳۰	-۰/۳۲	۱/۲۸
۴۹	۶/۱۵	۶/۱۲	۶/۰۷	-۰/۴۸	-۱/۳۰
۵۰	۵/۸۸	۵/۸۷	۵/۹۹	-۰/۱۷	۱/۸۷
۵۱	۶	۵/۹۴	۵/۹۷	-۱/۰۰	-۰/۵۰
۵۲	۶/۳۰	۶/۴۰	۶/۱۷	۱/۵۸	-۲/۰۶
۵۳	۵/۸۸	۵/۹۰	۵/۹۳	۰/۳۴	۰/۸۵



شکل ۳-۲۰: نمودار مقادیر پیش‌بینی شده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک توسط مدل شبکه عصبی مصنوعی بر حسب مقادیر تجربی



شکل ۳-۲۱: نمودار مقادیر خطای باقی‌مانده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک توسط مدل شبکه عصبی مصنوعی بر حسب مقادیر تجربی

### ۳-۷-۴ - مقایسه مدل‌های ارائه شده با استفاده از پارامترهای آماری

همان‌طور که قبلاً بیان شد، پارامترهای آماری مقادیر تعیین‌کننده‌ای برای ارزیابی مدل‌ها هستند. مطابق جدول ۳-۱۵ از این مقادیر برای ارزیابی و مقایسه مدل‌های ساخته شده استفاده گردید. این پارامترها برای سری ارزیابی، سری آزمون و کل داده‌ها به روش رد مرحله‌ای تک تک و مطابق با روابط توضیح داده شده در بخش ۲-۱۵-۱ محاسبه شده‌اند.

جدول ۳-۱۵: پارامترهای آماری محاسبه شده برای ارزیابی و مقایسه مدل‌های SR-ANN ، GC-ANN ، SR-RF و GC-RF

		SR-ANN	GC-ANN	SR-RF	GC-RF
R	سری ارزیابی	۰/۹۸۳	۰/۹۷۳	—	—
	سری تست	۰/۹۶۷	۰/۹۶۳	۰/۹۴۸	۰/۹۶۳
	کل داده‌ها	۰/۹۳۴	۰/۹۳۹	۰/۹۳۲	۰/۹۳۴
R <sup>2</sup>	سری ارزیابی	۰/۹۴۱	۰/۹۴۱	—	—
	سری تست	۰/۹۳۵	۰/۹۲۷	۰/۸۹۲	۰/۹۲۴
	کل داده‌ها	۰/۸۷۲	۰/۸۸۲	۰/۸۶۹	۰/۸۷۲
PRESS	سری ارزیابی	۰/۰۲۰	۰/۰۲۰	—	—
	سری تست	۰/۰۴۰	۰/۰۵۴	۰/۰۹۴	۰/۰۶۶
	کل داده‌ها	۰/۳۵۹	۰/۲۷۷	۰/۳۱۴	۰/۳۱۴
MSE	سری ارزیابی	۰/۰۰۲	۰/۰۰۳	—	—
	سری تست	۰/۰۰۵	۰/۰۰۷	۰/۰۱۷	۰/۰۰۷
	کل داده‌ها	۰/۰۰۷	۰/۰۰۵	۰/۰۰۶	۰/۰۰۶
MAE	سری ارزیابی	۰/۰۰۶	۰/۰۴۷	—	—
	سری تست	۰/۰۰۸	۰/۰۷۷	۰/۰۸۹	۰/۰۷۶
	کل داده‌ها	۰/۰۰۱	۰/۰۶۲	۰/۰۶۷	۰/۰۶۸
MRE	سری ارزیابی	۰/۰۶۳	۰/۷۹۸	—	—
	سری تست	۰/۵۰۰	۱/۱۳۵	۱/۵۱۰	۱/۲۹۱
	کل داده‌ها	۰/۰۳۸	۱/۰۴۱	۱/۱۳۱	۱/۱۴۲
SEP	سری ارزیابی	۰/۰۵۰	۰/۰۵۰	—	—
	سری تست	۰/۰۷۱	۰/۰۸۲	۰/۰۹۷	۰/۸۱۱
	کل داده‌ها	۰/۰۸۱	۰/۰۷۲	۰/۰۷۷	۰/۰۷۶۷
REP(%)	سری ارزیابی	۰/۸۴۹	۰/۸۵۲	—	—
	سری تست	۱/۱۹۰	۱/۳۸۶	۱/۶۴۸	۱/۳۸۶
	کل داده‌ها	۱/۳۴۹	۱/۲۲۰	۱/۲۹۵	۱/۲۹۹

### ۳-۸- بررسی ارتباط بین توصیف کننده‌های منتخب و فعالیت مورد نظر

در این قسمت با توجه به توصیف کننده‌های وارد شده در مدل‌های RF و ANN یک بررسی اجمالی روی اثرات این متغیرها بر فعالیت دارویی ترکیبات مورد مطالعه صورت خواهد گرفت. بهترین مدل انتخاب شده شامل ۸ توصیف کننده است که هر کدام بیانگر خصوصیت دو بعدی، هندسی و سه بعدی ترکیبات مورد بررسی می‌باشند. جدول ۳-۱۶ مقادیر اثر متوسط<sup>۱</sup> توصیف کننده‌های مدل برتر یعنی SR-ANN را نشان می‌دهد. اثر متوسط یک متغیر مستقل با استفاده از فرمول (۳-۱) به دست می‌آید:

(۳-۱)

$$T = \frac{\beta_x \sum X_n}{\sum pEC_{50}}$$

که در آن  $\beta_x$  ضریب متغیر مستقل X در مدل،  $X_n$  مقدار متغیر مستقل X مورد نظر برای ترکیب nام می‌باشد.

جدول ۳-۱۶: اثر متوسط توصیف کننده‌های به دست آمده توسط روش رگرسیون مرحله ای

Variable	B07NCl	Mor25m	MATS8m	Mor29v	RDF110u	TF..Cl	G1p	CIC5
Mean effect	-۸/۵۷۶	۶/۶۹۰	۳/۲۲۲	-۳/۵۵۰	-۳/۹۳۶	-۴/۵۶۹	-۳/۸۹۵	-۲/۹۳۲

### الف) توصیف کننده‌های 2D binary fingerprints

این توصیف کننده حضور یا عدم حضور N-Cl در فاصله توپولوژیکی ۷ را نشان می‌دهد. اثر انگشت باینری اندازه مولکول‌ها را به طور غیر مستقیم بررسی می‌کند. در حالت تجربی مولکول‌های بزرگتر، دارای ویژگی‌های متنوع‌تری هستند. در حالت تئوری فاصله توپولوژیکی مولکول‌های بزرگ ثابت است و اثر انگشت باینری حضور یا عدم حضور ویژگی‌های موثر در عملکرد مولکول را نسبت به حالتی که

<sup>۱</sup>- Mean effect

ویژگی‌های واقعی بررسی شده‌اند مورد بررسی قرار می‌دهد. توصیف کننده‌های اثر انگشتی اندازه مولکول را بررسی می‌کنند و همه این ویژگی‌های بررسی شده برای مولکول باید با اندازه مولکول در سطح اولیه همبسته باشد. که از میان توصیف کننده‌های اثر انگشت B07NCI در مدل برتر دیده می‌شود [۴۱].

### ب) توصیف کننده های 3D MORSE

توصیف کننده های 3D - Morse (نمایش سه بعدی ساختار مولکول براساس تفرق الکترون) از

طریق معادله تبدیلی که در پراش الکترون استفاده می‌شود، محاسبه می‌گردند:

$$I(s) = \sum_{i=2}^N \sum_{j=1}^{i-1} A_i A_j \frac{\sin(sr_{ij})}{sr_{ij}} \quad (2-3)$$

I شدت الکترون پراکنده شده،  $A_i$  و  $A_j$  خاصیت اتمی  $i$  و  $j$ ،  $s$  زاویه پراکندگی،  $r_{ij}$  فاصله بین

اتم‌های  $i$  و  $j$ ،  $N$  تعداد کل اتم‌ها را نشان می‌دهد. این روش باعث می‌شود که ساختار سه بعدی مولکول

به یک کد ثابت تبدیل شود. این توصیف کننده ها قادرند ارتباط بین ساختار سه بعدی ترکیبات آلی و

خصوصیات فیزیکی، شیمیایی و بیولوژیکی را فراهم سازند [۱۹]. برخی از این توصیف کننده های سه

بعدی مانند  $^1\text{Mor}25\text{m}$  و  $^2\text{Mor}29\text{v}$  در مدل برتر دیده می‌شوند، که به ترتیب توسط جرم اتمی و حجم

واندروالس وزن دار شده اند و این نشان می‌دهد که حجم کلی و جرم اتمی استخلاف، فاکتوری تعیین

کننده برای فعالیت دارویی آن ترکیب می‌باشد.

<sup>۱</sup> - 3D -MORSE-signal 25 weighted by atomic masses

<sup>۲</sup> - 3D -MORSE-signal 29 weighted by atomic van der waals volumes

### ب) توصیف کننده های Moran Autocorrelation

توصیف کننده های گروه Moran Autocorrelation توصیف کننده های دوبعدی هستند که از

طریق رابطه زیر محاسبه می شوند:

$$I(d) = \frac{\frac{1}{2} \sum_{i=1}^A \sum_{j=1}^A \delta_{ij} \cdot (w_i - \bar{w})(w_j - \bar{w})}{\frac{1}{A} \cdot \sum_{i=1}^A (w_i - \bar{w})^2} \quad (3-3)$$

که  $w$  یک ویژگی اتم،  $\bar{w}$  میانگین مقدار آن روی مولکول،  $A$  تعداد اتم ها و  $d$  فاصله ی توپولوژیکی است و  $\delta_{ij}$  نیز - که به تابع کرونکر<sup>۱</sup> معروف است - در حالتی که  $d_{ij} = d$  باشد یک است، در غیر اینصورت صفر است. مقدار این توصیف کننده که از جنس فاصله است، از  $+1$  تا  $-1$  متغیر است [۱۹ و ۴۲]. MATS8m<sup>۲</sup> از این گروه توصیف کننده هاست که در مدل برتر ظاهر شده است که با جرم اتمی وزن دار شده در مدل انتخاب شده و نشان می دهد که وزن کلی استخلاف، فاکتوری تعیین کننده برای فعالیت دارویی آن ترکیب می باشد. علامت مثبت این توصیف کننده بیانگر این است که که با افزایش مقادیر این توصیف کننده pEC50 نیز افزایش می یابد. جدول ۳-۱۷ چگونگی ارتباط این توصیف کننده با فعالیت دارویی را برای تعدادی از مولکول های سری داده ها نشان می دهد.

جدول ۳-۱۷: نمایش ارتباط فعالیت با مقدار توصیف کننده MATS8m برای بعضی از ترکیبات مورد مطالعه

شماره ترکیب	مقدار MATS8m	pEC50
۱۱	-۰/۲۳۸	۵/۶۹
۲۴	-۰/۱۵۰	۵/۸۵
۳۶	۰/۰۲۰	۶/۰۹
۴۹	۰/۰۳۹	۶/۱۵

<sup>۱</sup>- Kronecker

<sup>۲</sup>- Moran autocorrelation-lag8 weighted by atomic masses



### ت) توصیف کننده‌های گروه RDF<sup>۱</sup>

می‌توان تابع توزیع شعاعی یک ترکیب شامل N اتم را به احتمال یافتن یک اتم در یک حجم کروی با شعاع r تعریف کرد و فرم کلی این تابع به صورت زیر است:

$$g(r) = f \sum_{i=1}^{N-1} \sum_{j>i}^N A_i A_j e^{\beta(r-r_{ij})^2} \quad (4-3)$$

که در این رابطه N تعداد اتم‌ها،  $A_i$  و  $A_j$  خاصیت اتمی ( الکترونگاتیویته، جرم اتمی، قطبش پذیری اتمی و...)،  $r_{ij}$  شعاع بین اتم i و j، f فاکتور مقیاس،  $\beta$  فاکتور دما که باعث حرکت و جا به جایی اتم‌ها است، می‌باشند [۳۷]. از میان این توصیف کننده‌ها RDF<sup>۱</sup> در مدل برتر دیده می‌شود.

### ث) توصیف کننده‌های Topological [۱۹]

این توصیف کننده بر اساس نمایش گراف مولکول می‌باشند. در این گراف‌ها هر نقطه نشان‌دهنده یک اتم بوده و خطوط بین نقاط نیز نشان‌دهنده پیوند شیمیایی بین اتم‌ها می‌باشد. معمولاً در گراف‌های مولکولی اتم هیدروژن را نشان نمی‌دهند. این توصیف کننده اطلاعاتی راجع به ساختمان، اندازه، شکل، تقارن، شاخه‌دار شدن، نحوه‌ی اتصال اتم‌ها و نوع اتم‌های موجود در یک مولکول در اختیار ما قرار می‌دهند. محاسبه‌ی این توصیف کننده به سادگی از روی ساختمان دوبعدی مولکول‌ها امکان‌پذیر می‌باشد. از این گروه، توصیف کننده‌های CIC5<sup>۲</sup> و TF.Cl<sup>۳</sup> انتخاب شده است.

### ج) توصیف کننده‌های WHIM<sup>۴</sup>

این شاخص از مختصات کارتزین ساختار سه‌بعدی مولکول، با استفاده از صورتبندی با حداقل انرژی

---

<sup>۱</sup> - Radial Distribution Function

<sup>۲</sup> - Complementary information content

<sup>۳</sup> - Sum of topological distances between F..Cl

<sup>۴</sup> - Weighted Holistic Invariant Molecular descriptors

محاسبه می‌شود و شامل اطلاعاتی درباره‌ی اندازه، شکل، تقارن و توزیع اتمی ساختار سه‌بعدی مولکول می‌باشد. این توصیف کننده از رابطه (۵-۳) به دست می‌آید:

$$S_{jk} = \frac{\sum_{i=1}^A w_i (q_{ij} - \bar{q}_j) (q_{ik} - \bar{q}_k)}{\sum_{i=1}^A w_i} \quad (5-3)$$

که  $S_{jk}$  کوواریانس وزن دار شده بین کئوردینه  $j$  ام و  $k$  ام،  $A$  تعداد اتم‌ها،  $w_i$  وزن  $i$  امین اتم،  $q_{ij}$  و  $q_{ik}$  به ترتیب  $j$  امین و  $k$  امین کئوردینه‌های اتم  $i$  ام، و  $\bar{q}$  مقدار میانگین متناظر است. شش طرح وزن دار شدن پیشنهاد شده است که عبارتند از:

- ۱ - حالت بدون وزن ( $u$ ) ۲- وزن دار شده با جرم اتمی ۳- وزن دار شده با حجم و اندروالس ۴-
- وزن دار شده با الکترونگاتیویته ساندرسون ۵- وزن دار شده با حالت الکتروتوپولوژیکی کی-یر-هال<sup>۱</sup> ۶-
- وزن دار شده با قطبش‌پذیری [۱۹ و ۴۳]. از این دسته توصیف کننده  $G1p$  توسط مدل انتخاب گردید که به تأثیر صورتبندی دارو اشاره دارد.

### ۳-۹- بررسی میزان مشارکت توصیف‌کننده‌های منتخب در شبکه عصبی

میزان مشارکت توصیف‌کننده‌های منتخب توسط شبکه عصبی (SR-ANN) به صورت زیر تعیین شد:

- ۱- توصیف‌کننده مورد نظر به همراه اوزان مربوطه اش از شبکه بهینه شده حذف گردید.
- ۲- با استفاده از بقیه توصیف‌کننده‌ها مقدار متغیر وابسته ( $PEC_{50}$ ) برای هر ترکیب سری ارزیابی پیش‌بینی شد.
- ۳- میانگین خطای مطلق ترکیبات سری ارزیابی محاسبه گردید.
- ۴- مراحل ۱ تا ۳ برای دیگر توصیف‌کننده‌های منتخب نیز تکرار شد.
- ۵- سرانجام درصد مشارکت هر توصیف‌کننده توسط رابطه (۳-۶) برآورد شد [۴۴ و ۴۵].

<sup>۱</sup> - Kier-hall

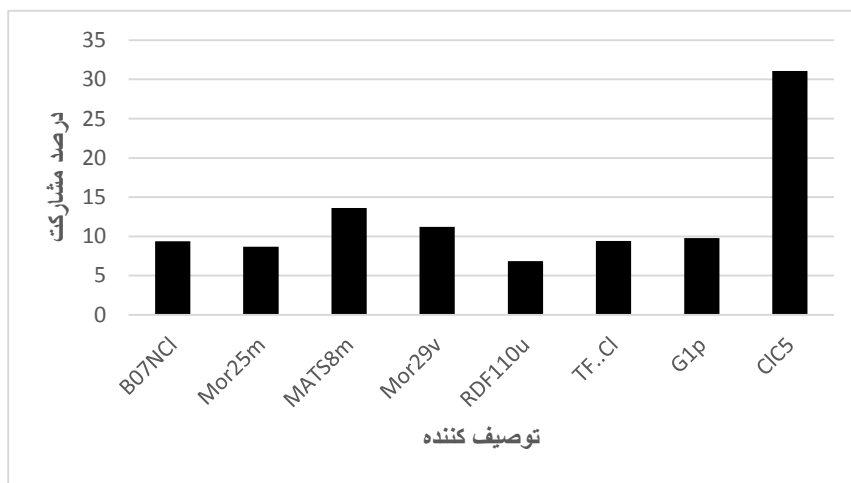
$$c_i = 100 \frac{\Delta m_i}{\sum_{i=1}^N \Delta m_i} \quad (6-3)$$

در این رابطه  $c_i$  درصد مشارکت توصیف‌کننده حذف شده  $i$ ام،  $N$  تعداد توصیف‌کننده‌های مدل

و  $\Delta m_i$  میانگین خطای مطلق سری ارزیابی در غیاب توصیف‌کننده  $i$ ام را نشان می‌دهد. بر این

اساس درصد مشارکت توصیف‌کننده‌های منتخب در ترکیبات مورد بررسی به صورت شکل ۳-۲۲ به

دست آمد.



شکل ۳-۲۲: مشارکت توصیف‌کننده‌ها در شبکه عصبی بهینه SR-ANN

بر اساس نمودار فوق توصیف‌کننده‌های ClC5، MATS8m و Mor29v دارای بیشترین اثر مشارکت می‌باشند و از میان این سه، توصیف‌کننده ClC5 دارای بیشترین اثر مشارکت می‌باشد. این توصیف‌کننده در برگیرنده اطلاعات دو بعدی مولکول است. بنابراین می‌توان گفت که فعالیت ضد ایدز این ترکیبات وابستگی قابل ملاحظه‌ای به توپولوژی ترکیب ضد ایدز دارد. توصیف‌کننده دیگری که درصد مشارکت بالایی دارد، توصیف‌کننده MATS8m از دسته توصیف‌کننده دو بعدی است که وزن کلی استخلاف را در میزان فعالیت ضد ایدز بیان می‌کند. همچنین توصیف‌کننده Mor29v در برگیرنده اطلاعات سه بعدی مولکول است که این اطلاعات به صورتبندی دارو ارتباط دارد پس می‌توان گفت که فعالیت بازداری این ترکیبات وابستگی قابل ملاحظه‌ای به تغییرات صورتبندی دارو دارد.

### ۳-۱۰- نتیجه‌گیری نهایی

با توجه به این که به ندرت مطالعات استفاده از توصیف کننده‌های سهم گروه برای پیش‌بینی خاصیت دارویی استفاده شده است، در این گزارش از روش‌های جنگل‌های تصادفی و شبکه عصبی مصنوعی با استفاده از توصیف کننده‌های سهم گروه برای مدل‌سازی استفاده گردید. این مدل‌ها می‌توانند به محققان در طراحی ترکیبات جدید از این مشتقات با استخلاف‌های مختلف کمک زیادی نمایند. توصیف کننده‌های سهم گروه از روی ساختار ترکیب به دست می‌آیند بنابراین ساده و قابل تفسیر می‌باشند. محور اصلی این پژوهش نیز ارائه توصیف کننده‌های ساده، در عین حال مرتبط با فعالیت مورد نظر است. با توجه به نتایج به دست آمده می‌توان گفت که فعالیت ضد ایدز این ترکیبات وابستگی قابل قبولی به ساختمان دو بعدی ترکیب دارد. در نتیجه ترکیبی می‌تواند فعالیت ضد ایدز بهتری داشته باشد که مقدار سه توصیف کننده مهم  $C1C5$ ،  $MATS8m$  و  $Mor29v$  برای آن زیاد باشد.

۴- فصل چهارم

# مطالعه ارتباط کمی ساختار-فعالیت

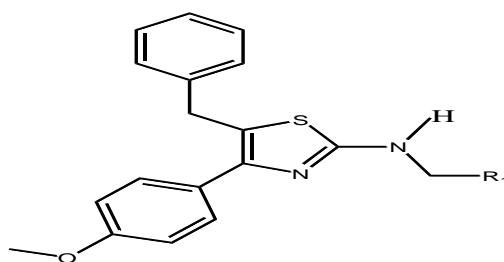
## مشتقات تiazol

## ۴-۱- مدل‌سازی فعالیت دارویی مشتقات تiazol

هدف در این فصل، یافتن یک مدل برای پیش‌بینی فعالیت ضد HIV مشتقاتی از ترکیبات تiazol است. بخش تجربی شامل معرفی سری داده‌ها، بهینه‌سازی ساختار مولکول‌ها و محاسبه توصیف‌کننده‌های مولکولی با استفاده از روش رگرسیون مرحله‌ای و همچنین مدل‌سازی توسط روش‌های جنگل‌های تصادفی و شبکه عصبی مصنوعی و ارزیابی مدل‌های برتر است.

## ۴-۲- سری داده‌ها

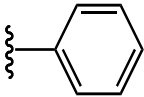
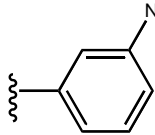
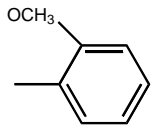

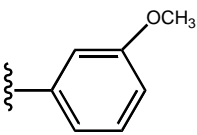
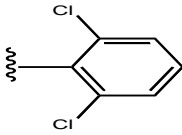
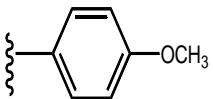
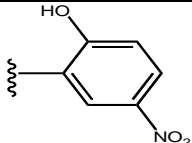
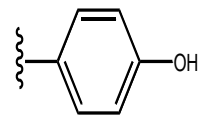
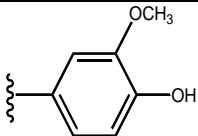
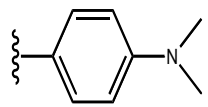
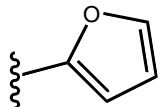
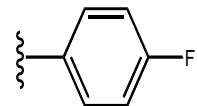
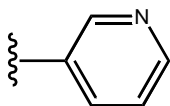
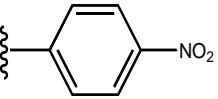
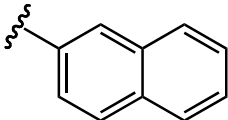
ترکیبات مورد استفاده در این تحقیق، داده‌های تجربی ۲۴ ترکیب از مشتقات تiazol می‌باشند که از مرجع [۴۶] گردآوری شده است. مقادیر منهای لگاریتم داده‌های  $IC_{50}$ <sup>۱</sup> (کمترین دوزی از دارو که آنزیم ۵۰٪ بازدارنده باشد) بر حسب مولار ( $-\log IC_{50} \times 10^{-6}$ ) به عنوان متغیر وابسته، در این مدل‌سازی مورد استفاده قرار گرفته است. دو ساختار اصلی برای این ترکیبات در شکل ۴-۱ و شکل ۴-۲ نشان داده شده است. جدول ۴-۱ و جدول ۴-۲ ساختار استخلاف‌های آن و مقادیر تجربی  $pIC_{50}$  بر حسب مولار که در شرایط یکسان اندازه‌گیری شده است را نشان می‌دهد.

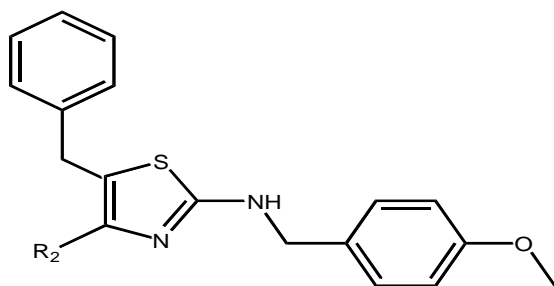


شکل ۴-۱: اسکلت اصلی شماره یک

<sup>۱</sup>- describes the dose that causes half-maximal inhibition by an antagonist

جدول ۴-۱ : ساختار و مقادیر تجربی  $pIC_{50}$  مشتقات

شماره	$R_1$	$pIC_{50}$	شماره	$R_1$	$pIC_{50}$
۱		۶/۱۹	۹		۵/۸۳
۲		۵/۱۶	۱۰		۶/۳۵
۳		۶/۰۹	۱۱		۵/۱۴
۴		۶/۵۱	۱۲		۵/۲۸
۵		۶/۱۹	۱۳		۶/۴۸
۶		۶/۰۶	۱۴		۶/۰۰
۷		۶/۱۸	۱۵		۵/۳۳
۸		۶/۱۰	۱۶		۶/۰۲



شکل ۴-۲: اسکلت اصلی شماره دو

جدول ۴-۲: ساختار و مقادیر تجربی pIC<sub>50</sub> مشتقات تiazول

شماره	R <sub>2</sub>	pIC <sub>50</sub>	شماره	R <sub>2</sub>	pIC <sub>50</sub>
۱۷		۷/۲۱	۲۱		۶/۸۵
۱۸		۵/۹۲	۲۲		۵/۲۱
۱۹		۷/۱۹	۲۳		۶/۱۵
۲۰		۷/۳۴	۲۴		۶/۷۴



### ۳-۴- بهینه‌سازی ساختمان هندسی مولکول‌ها و محاسبه توصیف‌کننده‌ها

در این مرحله، ساختار ترکیبات به وسیله نرم افزار Hyperchem رسم و برای بهینه کردن ساختار آن‌ها، روش AM<sub>1</sub> به کار برده شد. سپس بوسیله نرم افزار Dragon، ۳۲۲۴ توصیف‌کننده با استفاده از این ساختارهای بهینه شده محاسبه شد.

### ۴-۴- انتخاب بهترین توصیف‌کننده‌ها با روش رگرسیون

#### مرحله‌ای (SR)

برای انتخاب توصیف‌کننده‌هایی که به طور معناداری فعالیت دارویی را به ساختار مولکول ارتباط دهند، مراحل زیر انجام گرفت:

الف- توصیف‌کننده‌هایی که دارای مقادیر یکسان برای همه ترکیبات بودند، از فهرست توصیف‌کننده‌ها حذف شدند.

ب- همبستگی میان توصیف‌کننده‌ها محاسبه گردید و از میان جفت توصیف‌کننده‌هایی که ضریب همبستگی بزرگتر از ۰/۹ داشتند، آنکه که همبستگی کمتری با متغیر وابسته داشت، حذف شد.

ج- روش رگرسیون خطی چندگانه مرحله‌ای بر اساس انتخاب رو به جلو و حذف رو به عقب جهت غربال کردن بهترین توصیف‌کننده‌های مورد نیاز با استفاده از نرم افزار SPSS 18 انجام شد که در این مرحله ۱۰ توصیف‌کننده انتخاب شدند که بیشترین همبستگی با فعالیت ضد HIV ترکیبات مورد مطالعه را دارند. جدول ۳-۴ کل توصیف‌کننده‌های انتخاب شده و طبقه مربوط به آنها را نشان می‌دهد. همچنین ماتریس ضریب همبستگی این توصیف‌کننده‌ها که در جدول ۴-۴ ارائه شده، عدم همبستگی قابل ملاحظه بین توصیف‌کننده‌ها را نشان می‌دهد.

جدول ۴-۳: توصیف کننده‌های انتخاب شده توسط رگرسیون مرحله‌ای

No	Symbol	Class	Meaning
۱	P2u	WHIM descriptors	2nd component shape directional WHIM index/unweighted
۲	F07SCl	2D frequency fingerprints	Frequency of S - Cl at topological distance 07
۳	MATS5m	2D autocorrelation	Moaran autocorrelation –lag5/Weighted by atomic masses
۴	GATS6p	2D autocorrelation	Geary autocorrelation –lag 6/ Weighted by atomic polarizabilities
۵	SIC3	Information indices	Structural information content (neighborhood symmetry of 3-order)
۶	Mor30v	3D-MoRSE descriptors	3D-MoRSE- signal 30/Weighted by atomic van der waals volumes
۷	RDF055v	RDF descriptors	Radial Distribution Function -5.5/Weighted by atomic van der waals volumes
۸	G2p	WHIM descriptors	2st component symmetry directional WHIM index/ Weighted by atomic polarizabilities
۹	EEig08d	edge adjacency indices	Eigenvalue 08 from edge adj matrix weighted by dipol moments
۱۰	qpmax	Charge descriptors	Maximum positive charge

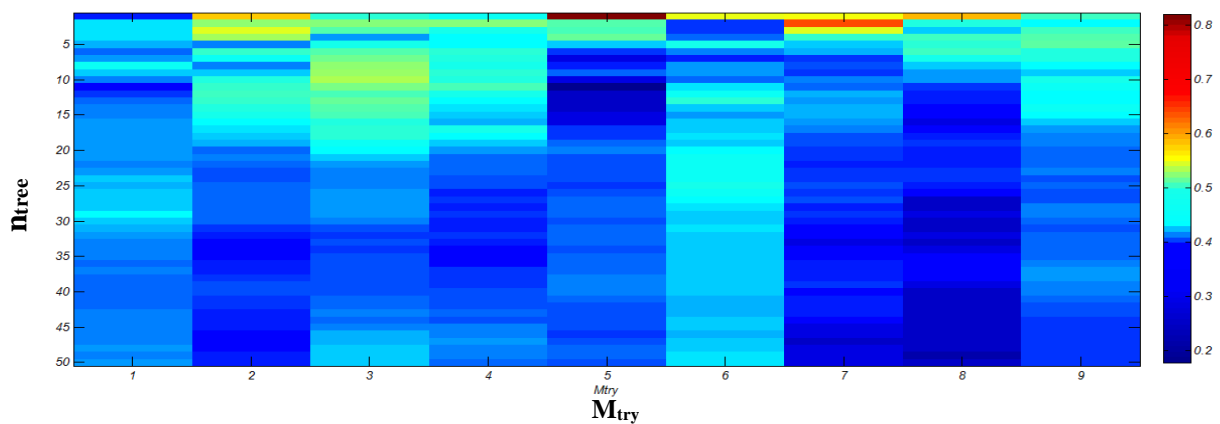
جدول ۴-۴: ماتریس همبستگی کل توصیف کننده‌های انتخاب شده توسط رگرسیون مرحله‌ای

	P2u	F07SC	MATS5m	GATS6p	SIC3	Mor30v	RDF055v	G2p	EEig08d	Qpmax
P2u	۱									
F07SC	۰/۲۵۰	۱								
MATS5m	۰/۴۱۱	۰/۴۴۱	۱							
GATS6p	۰/۰۵۳	۰/۱۸۲	۰/۲۵۱	۱						
SIC3	۰/۱۶۶	۰/۲۳۷	۰/۱۸۳	-۰/۲۹۹	۱					
Mor30v	۰/۲۵۲	-۰/۱۶۸	۰/۰۰۹	-۰/۰۶۹	-۰/۰۶۹	۱				
RDF055v	۰/۳۴۳	۰/۲۹۶	۰/۱۹۳	۰/۴۱۰	-۰/۲۰۹	-۰/۰۷۳	۱			
G2p	۰/۱۰۳	۰/۰۶۴	۰/۲۱۸	-۰/۲۷۸	-۰/۴۳۳	۰/۰۵۶	-۰/۱۸۰	۱		
EEig08d	۰/۲۲۴	۰/۱۴۷	-۰/۰۴۹	-۰/۱۸۱	۰/۲۲۵	۰/۳۰۸	۰/۴۱۵	۰/۰۶۶	۱	
Qpmax	۰/۰۷۷	-۰/۰۸۷	-۰/۱۶۱	-۰/۳۴۹	۰/۳۴۸	-۰/۰۱۲	-۰/۱۸۱	۰/۱۴۲	۰/۵۷۴	۱

#### ۴-۵- مدل سازی و بهینه سازی پارامترهای مؤثر بر روش جنگل‌های تصادفی

برای استفاده از قابلیت انتخاب توصیف کننده و انجام مدل سازی توسط جنگل‌های تصادفی، ابتدا باید پارامترهای مؤثر بهینه شوند. برای انجام این کار ابتدا ترکیبات به دو سری آموزش (۱۸ ترکیب) و سری آزمون (۶ ترکیب) تقسیم شدند. مقادیر مربوط به ۱۰ توصیف کننده به عنوان متغیر مستقل و مقادیر  $pIC_{50}$  به عنوان متغیر وابسته در نظر گرفته شدند. داده‌های سری OOB نیز برای بهینه سازی تعداد درختان ( $n_{tree}$ )، تعداد توصیف کننده‌های انتخاب شده در هر مرحله افزاز ( $M_{try}$ ) و تعداد مشاهدات باقی مانده در هر گره (Node size) در نظر گرفته

شدند. برای بهینه نمودن، تعداد درخت از ۱۰۰ تا ۵۰۰ با گام صد، تعداد توصیف‌گرهای انتخاب شده در هر مرحله افزایش ( $M_{try}$ ) از ۲ تا ۱۰ با گام یک و مقدار مشاهدات باقی‌مانده در هر گره (Node size) از ۱ تا ۱۰ تغییر داده شد و در هر مرحله مقدار خطای مربوط به مجموعه OOB محاسبه گردید. نتایج بهینه‌سازی پارامترهای  $M_{try}$  و  $n_{tree}$  در Node size بهینه در شکل ۳-۴ آمده است همانطور که مشاهده می‌شود نقاط آبی رنگ دارای خطای کمتری است. همچنین نتایج بهینه‌سازی در جدول ۴-۵ برای چندین حالت که کمترین MSE را دارد، آورده شده است.



شکل ۳-۴: بهینه‌سازی پارامترهای  $M_{try}$  و  $n_{tree}$  در Node size بهینه

جدول ۴-۵: کمترین مقادیر MSE همراه با  $M_{try}$  و  $n_{tree}$  متناظر با آنها

$n_{tree}$	$M_{try}$	Node Size	MSE OOB
۲۰۰	۸	۱	۰/۰۳۹۸
۱۰۰	۶	۲	۰/۰۴۳۴
۱۰۰	۱۰	۳	۰/۰۴۹۷
۲۰۰	۸	۴	۰/۰۵۲۸
۵۰۰	۱۰	۵	۰/۰۵۹۴
۲۰۰	۷	۶	۰/۰۵۸۸
۱۰۰	۱۰	۷	۰/۰۶۴۷
۱۰۰	۱۰	۸	۰/۰۷۳۵
۵۰۰	۱۰	۹	۰/۰۹۷۸
۳۰۰	۱۰	۱۰	۰/۱۲۷۱

با توجه به نتایج حاصل از جدول، تعداد درختان = ۲۰۰ و تعداد توصیف کننده‌های انتخاب شده در هر مرحله افزاز = ۸ و تعداد گره = ۱ به عنوان مقادیر بهینه انتخاب شدند.

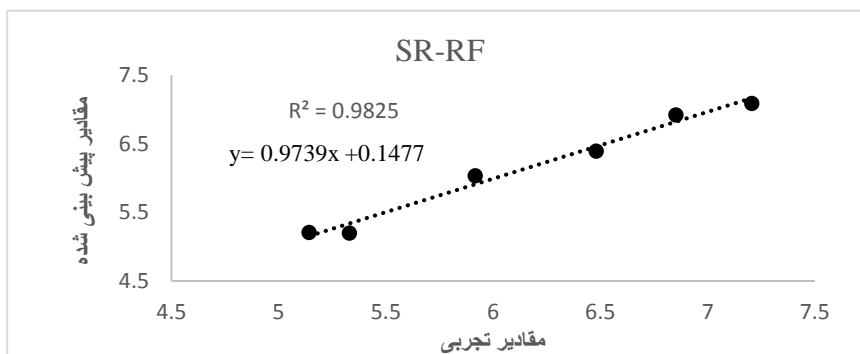
## ۴-۶- ارزیابی مدل جنگل‌های تصادفی

### ۴-۶-۱- ارزیابی مدل RF با استفاده از داده‌های سری آزمون

اهمیت مدل‌های بهینه وقتی مشخص می‌گردد که خواص مولکول‌هایی که در مدل‌سازی به کار نرفته‌اند را پیش بینی کند. بدین منظور مدل‌های منتخب، برای پیش‌بینی فعالیت داده‌های سری آزمون به کار برده شدند. جدول ۴-۶ نتایج پیش‌بینی حاصل از ارزیابی مدل جنگل‌های تصادفی با استفاده از داده‌های سری آزمون و شکل ۴-۴ نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی را برای داده‌های سری آزمون نشان می‌دهد.

جدول ۴-۶: نتایج حاصل از ارزیابی مدل جنگل تصادفی به دست آمده با استفاده از داده‌های سری آزمون

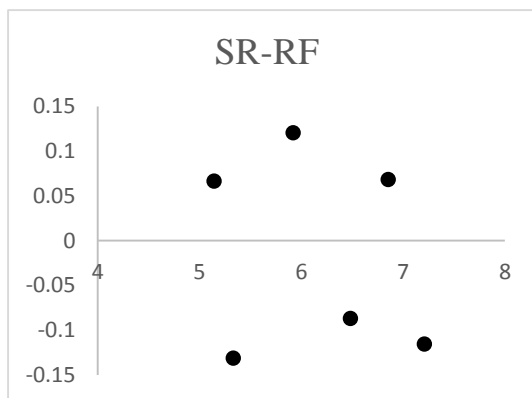
شماره ترکیب	مقدار تجربی	مقدار پیش‌بینی	درصد خطا
۱۱	۵/۱۴	۵/۲۱	۱/۳۶
۱۳	۶/۴۸	۶/۳۹	-۱/۳۸
۱۵	۵/۳۳	۵/۲۰	-۲/۴۳
۱۷	۷/۲۱	۷/۰۹	-۱/۶۶
۱۸	۵/۹۲	۶/۰۴	۲/۰۲
۲۱	۶/۸۵	۶/۹۲	۱/۰۲



شکل ۴-۴: نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی به وسیله داده‌های سری آزمون مدل جنگل‌های تصادفی

#### ۴-۶-۲- ارزیابی مدل RF با استفاده از نمودار خطای باقی مانده

اختلاف مقادیر پیش‌بینی شده و مقادیر تجربی، خطای باقیمانده نامیده می‌شود. توزیع متقارن داده‌ها حول محور افقی (خطای صفر) حاکی از عدم وجود خطای سیستماتیک است. نمودار خطای باقیمانده برحسب مقادیر تجربی، برای مدل‌های ذکر شده در شکل ۴-۵ نشان داده شده است.



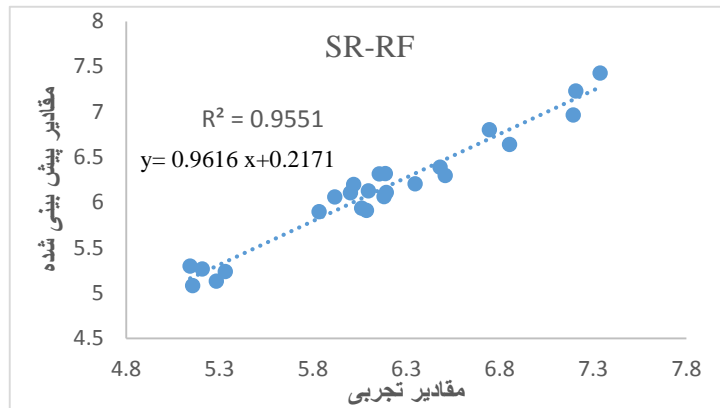
شکل ۴-۵: نمودار خطای باقیمانده برحسب مقادیر تجربی به وسیله داده‌های سری آزمون برای مدل جنگل‌های تصادفی

#### ۴-۶-۳- ارزیابی مدل RF توسط روش رد مرحله‌ای تک تک

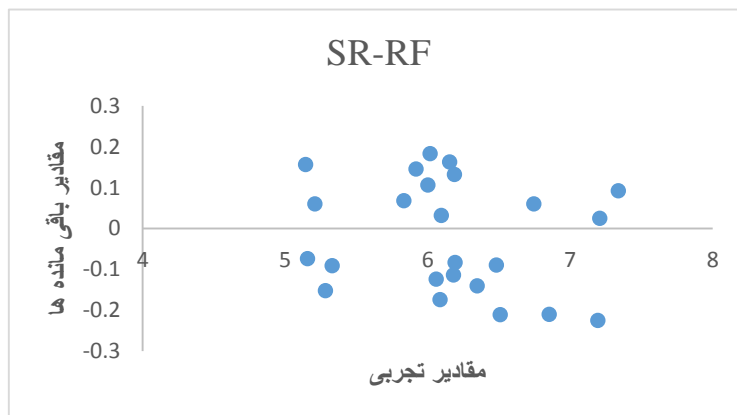
در این روش هر بار یک ترکیب به عنوان سری آزمون کنار گذاشته شد و در شرایط بهینه به دست آمده از بخش ۴-۵ با ترکیبات باقی مانده مدل‌سازی صورت گرفت. سپس مدل بدست آمده برای پیش‌بینی pIC50 ترکیب کنار گذاشته شده به کار گرفته شد و این فرآیند برای تمام اعضای سری داده‌ها تکرار شد. نتایج حاصل از این روش در جدول ۴-۷ نشان شده است. ضریب‌های تعیین مشاهده شده در شکل ۴-۶ نشان می‌دهند که مدل‌ها توانایی خوبی برای پیش‌بینی فعالیت ضد ایدز مشتقات تiazول دارد. همچنین نتایج به دست آمده شکل ۴-۷ که در آن‌ها مقادیر باقی مانده‌ها بر حسب مقادیر تجربی pIC50 ترکیبات مورد بحث ترسیم شده است دارای تقارن پراکندگی نقاط در دو طرف محور افقی است که نشان دهنده‌ی عدم وجود خطای سیستماتیک برای بیشتر ترکیبات است.

جدول ۴-۷: نتایج حاصل از رد مرحله‌ای تک‌تک برای مدل جنگل‌های تصادفی

شماره ترکیب	مقدار تجربی (pIC <sub>50</sub> )	مقدار پیش‌بینی	درصد خطا
۱	۶/۱۹	۶/۳۲	۲/۱۰
۲	۵/۱۶	۵/۰۸	-۱/۵۵
۳	۶/۰۹	۵/۹۱	-۲/۹۵
۴	۶/۵۱	۶/۳۰	-۳/۲۲
۵	۶/۱۹	۶/۱۱	-۱/۲۹
۶	۶/۰۶	۵/۹۴	-۲/۹۸
۷	۶/۱۸	۶/۰۷	-۱/۷۷
۸	۶/۱۰	۶/۱۳	۰/۴۹
۹	۵/۸۳	۵/۹۰	۱/۲۰
۱۰	۶/۳۵	۶/۲۱	-۲/۲۰
۱۱	۵/۱۴	۵/۳۰	۳/۱۱
۱۲	۵/۲۸	۵/۱۳	-۲/۸۴
۱۳	۶/۴۸	۶/۳۹	-۱/۳۸
۱۴	۶/۰۰	۶/۱۱	۱/۸۳
۱۵	۵/۳۳	۵/۲۴	-۱/۶۸
۱۶	۶/۰۲	۶/۲۰	۲/۹۹
۱۷	۷/۲۱	۷/۲۳	۰/۲۷
۱۸	۵/۹۲	۶/۰۶	۲/۳۶
۱۹	۷/۱۹	۶/۹۷	-۳/۰۵
۲۰	۷/۳۴	۷/۴۳	۱/۲۲
۲۱	۶/۸۵	۶/۶۴	-۳/۰۶
۲۲	۵/۲۱	۵/۲۷	۱/۱۵
۲۳	۶/۱۵	۶/۳۲	۲/۷۶
۲۴	۶/۷۴	۶/۸۰	۰/۸۹



شکل ۴-۶: نمودار مقادیر پیش‌بینی شده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک بر حسب مقادیر تجربی برای مدل جنگل‌های تصادفی



شکل ۴-۷: نمودار مقادیر خطای باقی‌مانده برای کل داده‌ها به روش رد مرحله‌ای تک‌تک بر حسب مقادیر تجربی برای مدل جنگل‌های تصادفی

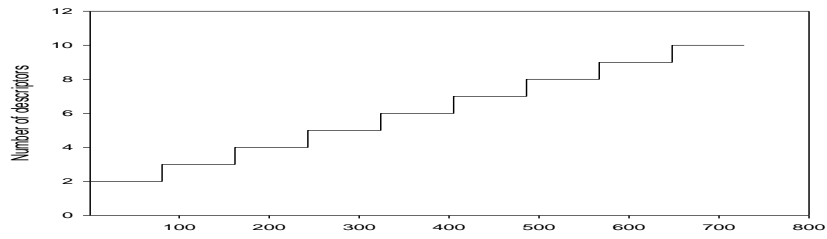


#### ۴-۷- مدل‌سازی و بهینه‌سازی پارامترهای موثر بر شبکه عصبی مصنوعی

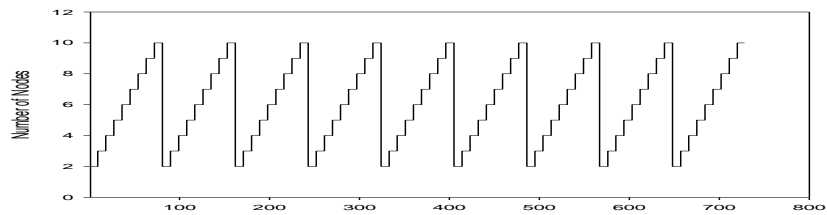
شبکه عصبی انتخاب شده در این مطالعه یک شبکه پیشخور با الگوریتم آموزشی پس انتشار می‌باشد. الگوریتم آن با برنامه متلب نوشته شده است. ابتدا مقادیر ۱۰ توصیف‌گر انتخاب شده با روش رگرسیون مرحله‌ای به عنوان متغیرهای ورودی و فعالیت ضد HIV متناظر آن‌ها به عنوان متغیر هدف در نظر گرفته شد تا پاسخ شبکه با آن‌ها سنجیده شود. برای به دست آوردن بهترین مدل و کم‌ترین خطا، پارامترهای مؤثر (تابع انتقال، تابع آموزش، تعداد متغیرهای ورودی شبکه، تعداد گره‌ها در لایه مخفی و تعداد دوره‌های آموزش) به‌طور همزمان بهینه‌سازی شدند.

در فرآیند بهینه‌سازی پارامترهای شبکه، سری داده‌ها به مجموعه آزمون (۶ ترکیب) و مجموعه آموزش (۱۸ ترکیب) تقسیم شدند. سری آموزش برای محاسبه گرادیان و به روز کردن اوزان و بایاس‌های شبکه با روش ارزیابی تقاطعی استفاده شد. از آنجایی که در بیشتر موارد به نظر می‌رسد یک لایه پنهان مناسب باشد [۴۰] در این تحقیق نیز از یک لایه پنهان استفاده شده است. برای یافتن مقدار بهینه پارامترهای مؤثر، مقادیر مختلفی از آنها در ساختار شبکه قرار داده شد. برای این منظور هر شبکه با تعداد ورودی ۲ تا ۱۰ و با دو الگوریتم آموزشی لونیبرگ-مارکوات (trainlm) و تنظیم بایزین (trainbr) و دو تابع انتقال لگاریتم سیگموئید و تانژانت سیگموئید و تعداد گره از ۲ تا ۱۰ و تعداد دوره‌های آموزش از ۱۰ تا ۵۰ با گام ۵، به‌طور همزمان آموزش داده شد.

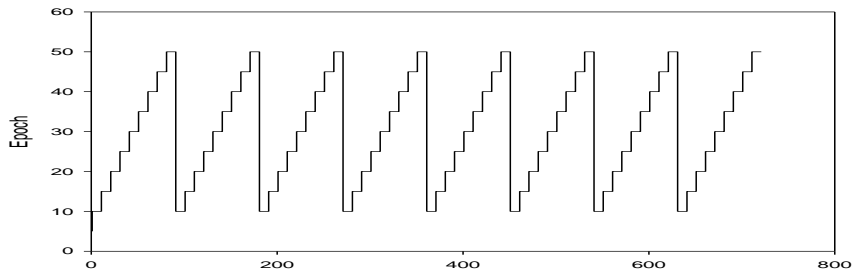
در روند بهینه‌سازی پارامترهای فوق، به حداقل رساندن میانگین مربع خطاهای حاصل از ارزیابی تقاطعی سری آموزش به عنوان معیار انتخاب شد. بخشی از روند تغییرات پارامترهای شبکه در حین بهینه‌سازی همزمان پارامترها به همراه مقادیر MSE بدست آمده، به صورت نموداری برحسب یک بردار مرجع فرضی در شکل ۴-۸ تا شکل ۴-۱۱ آمده است. نتایج بهینه‌سازی این پارامترها در جدول ۴-۸، گردآوری شده است.



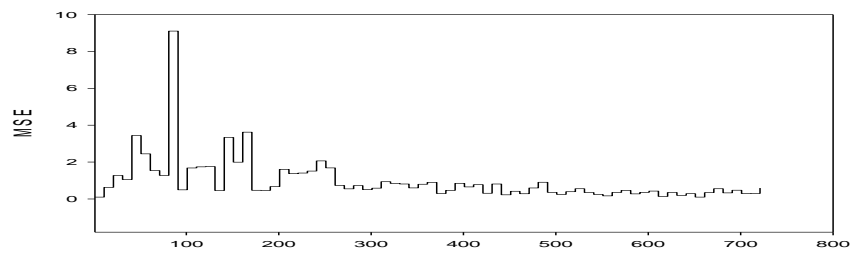
(الف)



(ب)



(ج)

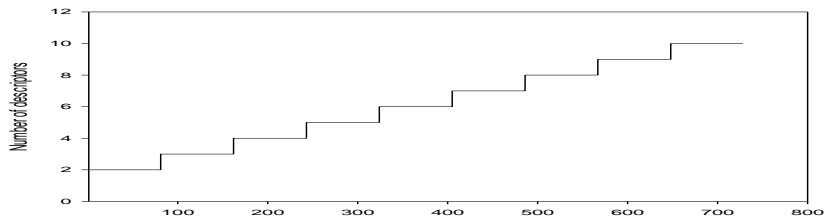


بردار مرجع

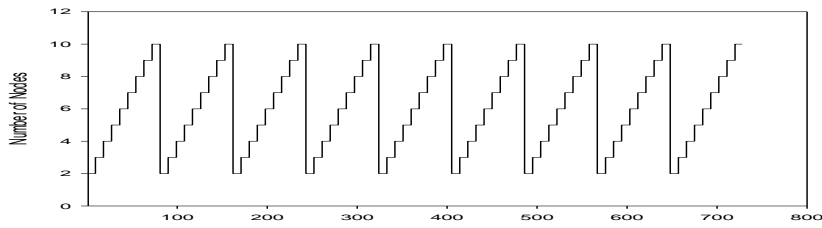
(د)

شکل ۴-۸: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دورهای آموزش و د) مقادیر MSE برای

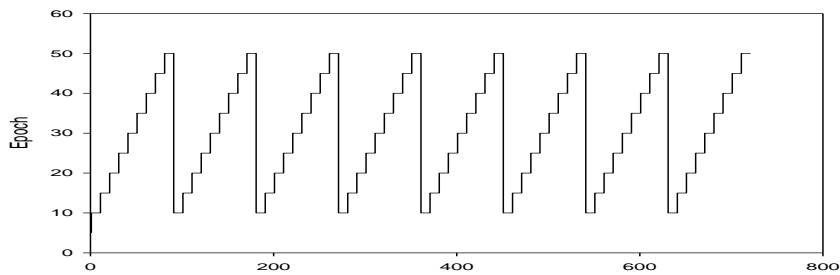
تابع آموزش لونبرگ-مارکوات به همراه تابع انتقال لگاریتم سیگموئید



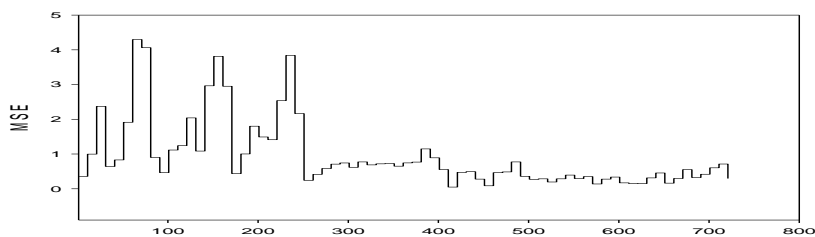
(الف)



(ب)



(ج)

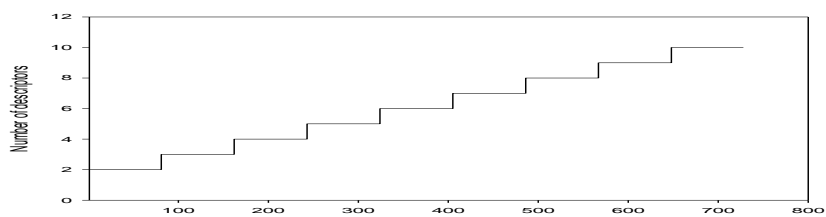


بردار مرجع

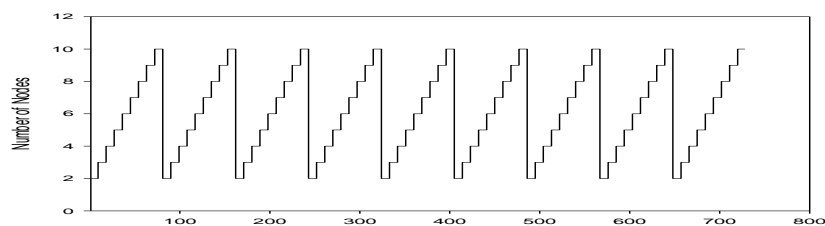
(د)

شکل ۴-۹: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای

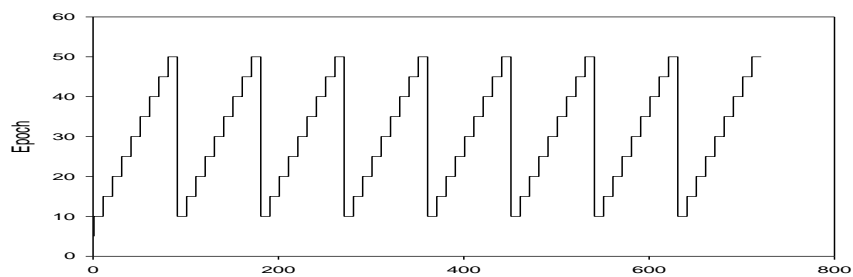
تابع آموزش لونبرگ-مارکوات به همراه تابع انتقال تانژانت سیگموئید



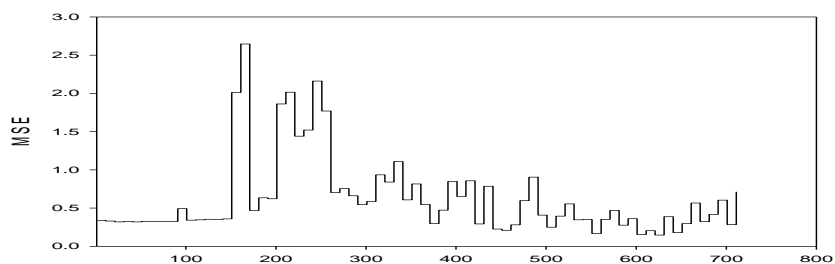
(الف)



(ب)



(ج)

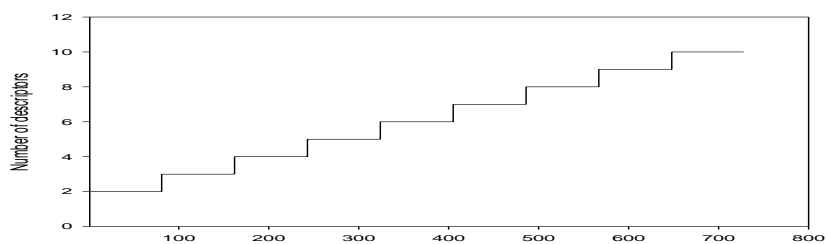


بردار مرجع

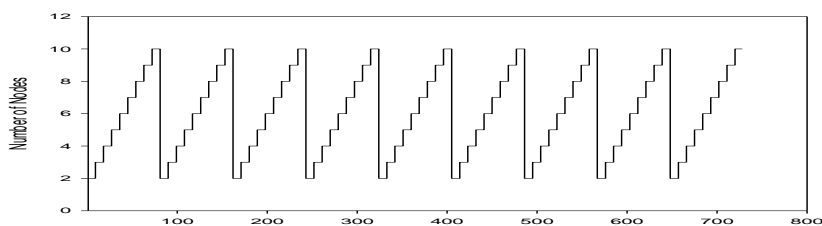
(د)

شکل ۴-۱۰: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دوره‌های آموزش و د) مقادیر MSE برای

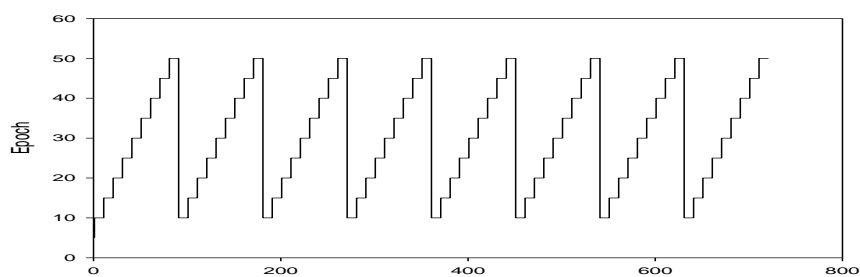
تابع آموزش تنظیم بایزین به همراه تابع انتقال لگاریتم سیگموئید



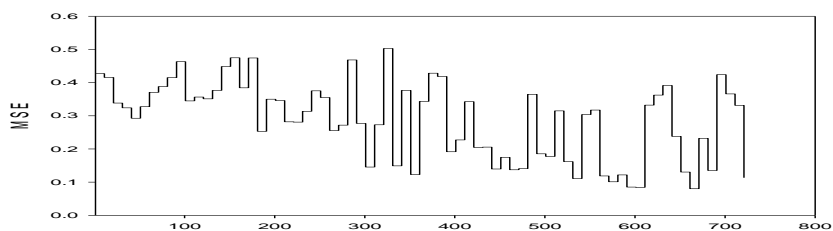
(الف)



(ب)



(ج)



(د)

بردار مرجع

شکل ۴-۱۱: نمودارهای الف) تعداد توصیف‌کننده، ب) تعداد گره‌ها، ج) تعداد دورهای آموزش و د) مقادیر MSE برای

تابع آموزش تنظیم بایزین به همراه تابع انتقال تانزانت سیگموئید

جدول ۴-۸: توابع و پارامترهای شبکه‌های بهینه SR-ANN

MSE	تعداد دور آموزش	تعداد گره	تابع انتقال	تابع آموزش	تعداد توصیف کننده
۰/۰۹۹۵	۱۰	۲	لگاریتم-سیگموئید	لونبرگ-مارکوات	۱۰
۰/۰۴۶۴	۴۰	۲	تانژانت-سیگموئید	لونبرگ-مارکوات	۷
۰/۰۸۳۶	۴۰	۴	لگاریتم-سیگموئید	تنظیم بایزین	۹
۰/۰۶۹۰	۱۵	۱۰	تانژانت-سیگموئید	تنظیم بایزین	۱۰

با توجه به نتایج بدست آمده، الگوریتم آموزشی لونبرگ-مارکوات و تابع انتقال تانژانت سیگموئید، ۷ توصیف کننده ورودی و ۲ نرون در لایه مخفی با تعداد دور آموزشی ۴۰ کمترین MSE را نشان می‌دهد. بنابراین این شبکه برای مدل‌سازی در نظر گرفته شد.

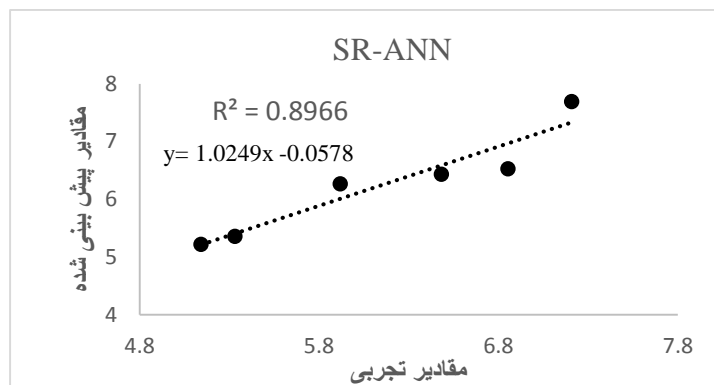
#### ۸-۴-۸- ارزیابی مدل شبکه عصبی مصنوعی

#### ۴-۸-۱- ارزیابی مدل ANN با استفاده از داده‌های سری آزمون

جدول ۴-۹ نتایج پیش‌بینی حاصل از ارزیابی مدل‌های شبکه عصبی مصنوعی با استفاده از داده‌های سری آزمون و شکل ۴-۱۲ نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی را برای داده‌های سری آزمون نشان می‌دهد.

جدول ۴-۹: نتایج حاصل از ارزیابی مدل شبکه عصبی مصنوعی با استفاده از داده‌های سری آزمون

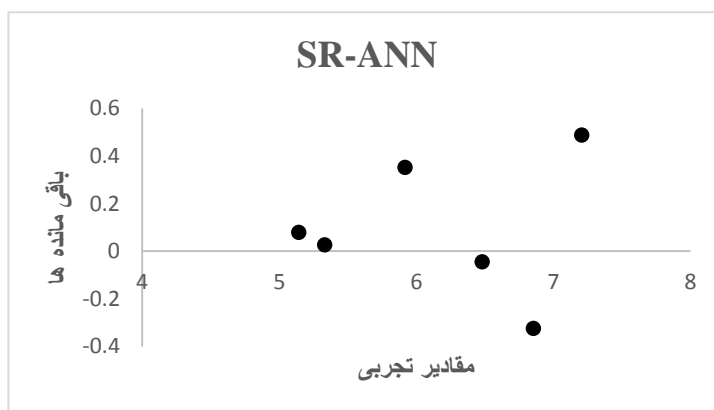
شماره ترکیب	مقدار تجربی (pIC <sub>50</sub> )	مقدار پیش‌بینی	درصد خطا
۱۱	۵/۱۴	۵/۲۲	۱/۵۵
۱۳	۶/۴۸	۶/۴۴	-۰/۶۱
۱۵	۵/۳۳	۵/۳۶	۰/۵۶
۱۷	۷/۲۱	۷/۶۹	۶/۶۵
۱۸	۵/۹۲	۶/۲۷	۵/۹۱
۲۱	۶/۸۵	۶/۵۳	-۴/۶۷



شکل ۴-۱۲: نمودار تغییرات مقادیر پیش بینی شده در مقابل مقادیر تجربی به وسیله داده‌های سری آزمون برای مدل شبکه عصبی مصنوعی

#### ۴-۸-۲- ارزیابی مدل ANN با استفاده از نمودار خطای باقی مانده

اختلاف مقادیر پیش‌بینی شده و مقادیر تجربی، خطای باقیمانده نامیده می‌شود. توزیع متقارن داده‌ها حول محور افقی (خطای صفر) حاکی از عدم وجود خطای سیستماتیک است. نمودار خطای باقیمانده برحسب مقادیر تجربی، برای مدل‌های ذکر شده در شکل ۴-۱۳ نشان داده شده است.



شکل ۴-۱۳: نمودار خطای باقیمانده برحسب مقادیر تجربی برای داده‌های سری آزمون برای مدل شبکه عصبی مصنوعی

#### ۴-۸-۳- ارزیابی مدل ANN توسط روش رد مرحله‌ای تک تک

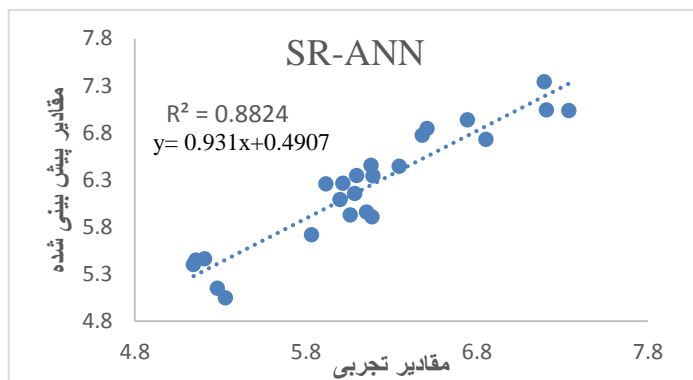
در این روش هر بار یک ترکیب به عنوان سری آزمون کنار گذاشته شد و در شرایط بهینه به دست آمده از بخش ۴-۷ با ترکیبات باقی مانده مدل‌سازی صورت گرفت. سپس مدل بدست آمده برای

پیش‌بینی  $pIC_{50}$  ترکیب کنار گذاشته شده به کار گرفته شد و این فرآیند برای تمام اعضای سری داده‌ها تکرار شد. نتایج حاصل از این روش در جدول ۴-۱۰ نشان شده است. ضریب تعیین مشاهده در شکل ۴-۱۴ نشان می‌دهد که مدل توانایی خوبی برای پیش‌بینی فعالیت ضد ایدز مشتقات تiazol دارد. همچنین نتایج به دست آمده شکل ۴-۱۵ که در آن مقادیر باقی‌مانده‌ها بر حسب مقادیر تجربی  $pIC_{50}$  ترکیبات مورد بحث ترسیم شده است دارای تقارن پراکندگی نقاط در دو طرف محور افقی است که نشان دهنده‌ی عدم وجود خطای سیستماتیک برای بیشتر ترکیبات است.

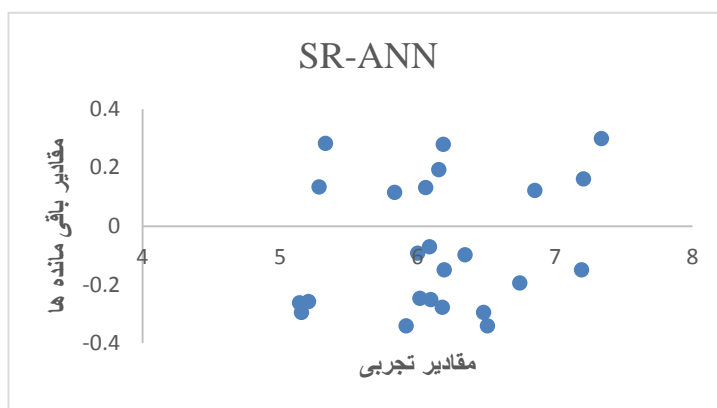


جدول ۴-۱۰: نتایج حاصل از رد مرحله‌ای تک‌تک برای مدل شبکه عصبی

شماره ترکیب	مقدار تجربی (pIC50)	مقدار پیش‌بینی	درصد خطا
۱	۶/۱۹	۵/۹۱	-۴/۵۲
۲	۵/۱۶	۵/۴۵	۵/۶۲
۳	۶/۰۹	۶/۱۶	۱/۱۴
۴	۶/۵۱	۶/۸۵	۵/۲۲
۵	۶/۱۹	۶/۳۴	۲/۴۲
۶	۶/۰۶	۵/۹۳	-۲/۱۴
۷	۶/۱۸	۶/۴۶	۴/۵۳
۸	۶/۱۰	۶/۳۵	۴/۰۹
۹	۵/۸۳	۵/۷۲	-۱/۸۸
۱۰	۶/۳۵	۶/۴۴	۱/۴۱
۱۱	۵/۱۴	۵/۴۰	۵/۰۵
۱۲	۵/۲۸	۵/۱۵	-۲/۴۶
۱۳	۶/۴۸	۶/۷۸	۴/۶۲
۱۴	۶/۰۰	۶/۰۹	۱/۵۰
۱۵	۵/۳۳	۵/۰۵	-۵/۲۵
۱۶	۶/۰۲	۶/۲۶	۳/۹۸
۱۷	۷/۲۱	۷/۰۵	-۲/۲۱
۱۸	۵/۹۲	۶/۲۶	۵/۷۴
۱۹	۷/۱۹	۷/۳۴	۲/۰۸
۲۰	۷/۳۴	۷/۰۴	-۴/۰۸
۲۱	۶/۸۵	۶/۷۳	-۱/۷۵
۲۲	۵/۲۱	۵/۴۶	۴/۷۹
۲۳	۶/۱۵	۵/۹۶	-۳/۰۸
۲۴	۶/۷۴	۶/۹۴	۲/۹۶



شکل ۴-۱۴: نمودار مقادیر پیش‌بینی شده برای کل داده‌ها به وسیله روش رد مرحله‌ای تک‌تک بر حسب مقادیر تجربی برای مدل شبکه عصبی مصنوعی



شکل ۴-۱۵: نمودار مقادیر خطای باقی‌مانده برای کل داده‌ها به وسیله روش رد مرحله‌ای تک‌تک بر حسب مقادیر تجربی مدل شبکه عصبی مصنوعی

#### ۴-۸-۴- مقایسه مدل‌ها با استفاده از پارامترهای آماری

همان‌طور که قبلاً بیان شد، پارامترهای آماری مقادیر تعیین‌کننده‌ای برای ارزیابی مدل‌ها هستند. مطابق جدول ۴-۱۱ در این مرحله از این مقادیر برای ارزیابی مدل‌های ساخته شده استفاده گردید. با توجه به نتایج حاصل از این جدول می‌توان گفت که مدل جنگل‌های تصادفی ارائه شده نسبت به شبکه عصبی مصنوعی برای مشتقات تیزول دارای برتری است.

جدول ۴-۱۱: پارامترهای آماری برای مدل‌های طراحی شده برتر توسط RF و ANN

پارامتر	سری تست (N=10)		کل داده‌ها (N=24)	
	RF	ANN	RF	ANN
<b>R</b>	۰/۹۹۱	۰/۹۴۷	۰/۹۷۷	۰/۹۳۹
<b>R<sup>2</sup></b>	۰/۹۸۲	۰/۸۶۲	۰/۹۵۴	۰/۸۶۸۲
<b>PRESS</b>	۰/۰۶۲	۰/۴۷۶	۰/۴۲۶	۱/۲۲۱۳
<b>MSE</b>	۰/۰۱۰	۰/۰۷۹۴	۰/۰۱۸	۰/۰۵۱
<b>MAE</b>	۰/۰۹۸	۰/۲۱۹۰	۰/۱۲۱	۰/۲۱۰
<b>MRE</b>	۱/۶۲۳	۳/۳۶۲	۱/۹۸۲	۳/۴۷۴
<b>SEP</b>	۰/۲۴۸	۰/۲۸۲	۰/۱۳۳۱	۰/۲۲۶
<b>REP(%)</b>	۴/۰۴۶	۴/۵۰۸	۲/۱۷۴	۳/۶۳۱

## ۹-۴- بررسی ارتباط بین توصیف کننده‌های منتخب و خاصیت موردنظر

در این بخش به بررسی ارتباط بین توصیف‌گرهای منتخب در مدل‌سازی جنگل‌های تصادفی که مدل برتر است و در بخش ۴-۷ محاسبه شدند، می‌پردازیم و سعی می‌کنیم تا حد امکان توضیحی برای ارتباط بین این توصیف‌گرها و فعالیت مورد نظر بیابیم. جدول ۴-۱۲ مقادیر اثر متوسط<sup>۱</sup> (که روش محاسبه آن در بخش (۳-۸) بیان شد) توصیف‌کننده‌های مدل برتر یعنی SR-RF را نشان می‌دهد.

جدول ۴-۱۲: اثر متوسط توصیف‌کننده‌های به دست آمده توسط روش رگرسیون مرحله ای

Variable	P2u	F07SCI	MATS5m	GATS6p	SIC3	Mor30v	RDF055v	G2p
Mean effect	-۳/۵۸۰	-۷/۹۹۹	۱۵/۱۶۹	-۱۰/۹۴۰	-۷/۰۴۳	۶/۵۴۳	-۵/۳۴۹	-۴/۹۰۶

### الف) توصیف کننده‌های WHIM<sup>۲</sup>

همان‌طور که در بخش ج-۸-۳ اشاره شد، این شاخص از مختصات کارتیزین ساختار سه‌بعدی مولکول، با استفاده از صورتبندی با حداقل انرژی محاسبه می‌شود و شامل اطلاعاتی درباره‌ی اندازه، شکل، تقارن و توزیع اتمی ساختار سه‌بعدی مولکول می‌باشد. از این نوع توصیف‌گرها P2u, G2p توسط مدل انتخاب گردید که به تأثیر صورتبندی دارو اشاره دارد.

### ب) توصیف کننده‌های 2D frequency fingerprints

همان‌طور که در بخش الف-۸-۳ در مورد توصیف کننده‌های اثر انگشت اشاره شد، این توصیف کننده‌ها اندازه مولکول‌ها را بررسی می‌کند. توصیف کننده اثر انگشت فرکانس که به بررسی مجموع

<sup>۱</sup>- Mean effect

<sup>۲</sup>- Weighted Holistic Invariant Molecular descriptors

فرکانس دو اتم در فاصله هندسی آن اتم‌ها می‌پردازد، اندازه اتم‌ها را به‌طور مستقیم بررسی می‌کند. که از میان توصیف‌کننده‌های اثر انگشت F07SCI در مدل برتر دیده می‌شود که فرکانس S-Cl را در فاصله هندسی ۷ نشان می‌دهد [۴۷].

### پ) توصیف‌کننده MATS5m

همان‌طور که در بخش پ-۸-۳ اشاره شد، توصیف‌کننده‌های گروه Moran Autocorrelation توصیف‌کننده‌های دوبعدی هستند MATS5m<sup>۱</sup> از این گروه توصیف‌کننده‌ها ست که در مدل برتر ظاهر شده است.

### ت) توصیف‌کننده GATS6p

جزء (2D-autocorrelation) توصیف‌کننده‌های دوبعدی خودارتباطی است که از گراف مولکولی و از طریق محاسبه مجموع اوزان اتم‌های انتهایی کل مسیرها با طول مسیر مورد نظر (lag)، به دست می‌آیند. نوعی از این توصیف‌کننده‌ها گروه Geary Autocorrelation است که ضریب گری نام دارد و بدین صورت محاسبه می‌شود:

$$C(d) = \frac{\frac{1}{2\Delta} \sum_{i=1}^A \sum_{j=1}^A \delta_{ij} \cdot (w_i - w_j)^2}{\frac{1}{A-1} \cdot \sum_{i=1}^A (w_i - \bar{w})^2} \quad (1-4)$$

که  $w$  یک ویژگی اتم،  $\bar{w}$  میانگین مقدار آن ویژگی روی مولکول،  $A$  تعداد اتم‌ها و  $d$  فاصله‌ی توپولوژیکی است و  $\delta_{ij}$  نیز - که به تابع کرونیگر<sup>۲</sup> معروف است - در حالتی که  $d_{ij} = d$  باشد یک است، و در غیر اینصورت صفر است و  $\Delta$  هم مجموعه‌ی  $\delta$  هاست. مقدار این توصیف‌کننده که از جنس فاصله

<sup>۱</sup> - Moran autocorrelation-lag15 weighted by atomic polarizabilities

<sup>۲</sup> - Kronecker

است، از صفر تا بی نهایت متغیر است [۱۹] از میان این توصیف کننده‌ها GATS6p و نقش قطبش‌پذیری اتمی ترکیبات را در جهت‌گیری به سوی مولکول هدف نشان می‌دهد. علامت منفی این توصیف کننده بیانگر این است که مقادیر بالاتر از این توصیف کننده برای فعالیت ضد HIV مطلوب نیست. جدول ۴-۱۳ چگونگی ارتباط این توصیف کننده با فعالیت دارویی را برای تعدادی از مولکول‌های سری داده‌ها نشان می‌دهد.

جدول ۴-۱۳: نمایش ارتباط فعالیت با مقدار GATS6p

شماره ترکیب	مقدار GATS6p	pIC <sub>50</sub>
۱۱	۱/۳۱۰	۵/۱۴
۱۵	۱/۱۶۶	۵/۳۳
۱۸	۰/۸۶۸	۵/۹۲

### ث) توصیف کننده Topological

همان‌طور که در بخش ث-۸-۳ بیان شد این توصیف کننده بر اساس نمایش گراف مولکول می‌باشند. در این گراف‌ها هر نقطه نشان‌دهنده یک اتم بوده و خطوط بین نقاط نیز نشان‌دهنده پیوند شیمیایی بین اتم‌ها می‌باشد. معمولاً در گراف‌های مولکولی اتم هیدروژن را نشان نمی‌دهند. این توصیف کننده اطلاعاتی راجع به ساختمان، اندازه، شکل، تقارن، شاخه‌دار شدن، نحوه‌ی اتصال اتم‌ها و نوع اتم‌های موجود در یک مولکول در اختیار ما قرار می‌دهند. محاسبه‌ی این توصیف کننده به سادگی از روی ساختمان دوبعدی مولکول‌ها امکان‌پذیر می‌باشد. از میان این توصیف کننده‌ها SIC<sup>۱</sup> در مدل برتر دیده می‌شود.

<sup>۱</sup> Structural information content

#### د) توصیف کننده‌های 3D – MORSE<sup>۱</sup>

همان‌طور که در بخش ب-۸-۳ بیان شد توصیف کننده‌های 3D - Morse ساختار سه بعدی مولکول را براساس تفرق الکترون نشان می‌دهند. از بین این توصیف کننده‌ها Mor30v<sup>۲</sup> در مدل برتر دیده می‌شود.

#### ه- توصیف کننده‌های RDF<sup>۳</sup>

همان‌طور که در بخش ت-۸-۳ بیان شد می‌توان تابع توزیع شعاعی یک ترکیب شامل N اتم را به احتمال یافتن یک اتم در یک حجم کروی با شعاع r تعریف کرد [۴۸]. اهمیت این توصیف کننده‌ها به دلیل اختلاف در توزیع اتم‌ها در مولکول‌ها و اثر این توزیع بر روی فعالیت مهارکنندگی ترکیبات می‌باشد. از میان این توصیف کننده‌ها RDF055v در مدل برتر دیده می‌شود.

---

<sup>۱</sup> - 3D- Molecular Representation of Structure based on Electron diffraction

<sup>۲</sup> - 3D –MORSE-signal 30 volumes

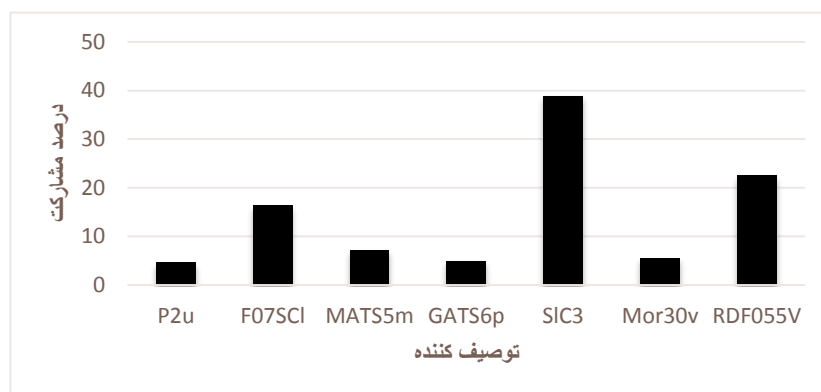
<sup>۳</sup> - Radial Distribution Function

## ۱۰-۴- بررسی میزان مشارکت توصیف کننده‌های منتخب در شبکه

### عصبی

بر اساس روشی که در بخش (۳-۹) توضیح داده شد، میزان مشارکت توصیف کننده‌های منتخب

در پیش‌بینی فعالیت دارویی موردنظر برآورد شد که نتایج آن در شکل ۴-۱۶ ارائه شده است.



شکل ۴-۱۶: مشارکت توصیف کننده‌ها در شبکه عصبی بهینه SR-ANN

بر اساس نمودار فوق توصیف کننده‌های SCI3، RDF055v و F07SCI دارای بیشترین اثر مشارکت می‌باشند. طبق شکل فوق توصیف کننده SCI3 دارای بیشترین اثر مشارکت می‌باشد. این توصیف کننده در برگیرنده اطلاعات دو بعدی مولکول است. توصیف کننده دیگری که درصد مشارکت بالایی دارد، توصیف کننده RDF055v که اهمیت این توصیف کننده به دلیل اختلاف در توزیع اتم‌ها در مولکول‌ها و اثر این توزیع بر روی فعالیت مهارکنندگی ترکیبات می‌باشد. همچنین توصیف کننده F07SCI بیان کننده اندازه مولکول به‌طور مستقیم است و اثر انگشت فرکانس را در فاصله هندسی بیان می‌کند.



#### ۴-۱۱- نتیجه‌گیری نهایی

در این تحقیق نتایج نشان می‌دهند که مدل‌های غیرخطی توسعه یافته می‌توانند به عنوان روش‌های موفق برای مدل‌سازی و پیش‌بینی فعالیت ضد ایدز ترکیبات مورد مطالعه به کار روند و این مدل‌ها می‌توانند به محققان در طراحی داروهایی با فعالیت ضد ایدز قویتر کمک کنند. بنابراین بر اساس این مدل‌ها می‌توان فعالیت ضد HIV را برای همین طبقه از ترکیبات که هنوز سنتز نشده‌اند پیش‌بینی کرد. در واقع این روش به طراحان دارو این امکان را می‌دهد که قبل از اینکه بخواهند یک ترکیب دارویی را سنتز کنند، ابتدا با استفاده از مدل معتبر خاصیت دارویی آن را پیش‌بینی کنند و در صورتی که دارای فعالیت دارویی مناسبی باشند، ترکیب سنتز شود و بدین ترتیب با صرف هزینه و وقت کمتر در مسیر تهیه داروهایی با اثر بخشی بهتر گام بردارند.

## پیشنهادات

- ✓ روش جنگل‌های تصادفی را می‌توان برای پیش‌بینی فعالیت سایر ترکیبات دارویی استفاده نمود.
- ✓ پژوهشگران می‌توانند با الگویی از روش سهم گروه مورد استفاده در این تحقیق، به بررسی روابط ساختار - فعالیت در گروه‌های دیگر مولکولی بپردازند.
- ✓ می‌توان از روش غیر خطی دیگر مانند SVM<sup>۱</sup> استفاده کرد.
- ✓ این تحقیقات را می‌توان با پژوهش روی خواص دارویی ترکیبات جدید ادامه داد تا نهایتاً منجر به ساخت داروهای جدیدتر با اثربخشی بهتر گردد.

---

<sup>۱</sup> - Supported Vector Machine

## فهرست منابع

- [1]-<http://www.soton.ac.uk/~ceb/teaching/2005/206-8.htm>.
- [2] -موسوی جاهد، " ویروس شناسی پزشکی"، انتشارات دانشگاه علوم پزشکی و پیراپزشکی شهید بهشتی، ۱۳۷۸.
- [3]-Debyser Z., Pauwels R., Andries K., Desmyter J., Kukla M., Janssen P. A., & An antiviral target on reverse transcriptase of human "De Clercq E. (1991), immunodeficiency virus type 1 revealed by tetrahydroimidazo-[4, 5, 1-jk][1, 4] , Proceedings of the National "benzodiazepin-2 (1H)-one and-thione derivatives Academy of Sciences, 88(4), 1451-1455.
- [4]- Pillai A. D., Rani S., Rathod P. D., Xavier F. P., Vasu K. K., Padh H., & QSAR studies on some thiophene analogs as anti-"Sudarsanam V. (2005), inflammatory agents: enhancement of activity by electronic parameters and its , Bioorganic & medicinal chemistry, 13(4), "utilization for chemical lead optimization 1275-1283.
- [5]- Jalali-Heravi M.; Parastar F., (2000), "Use of artificial neural network in a QSAR study of anti-HIV activity for a large group of HEPT derivatives", J. Chem. Inf. Comput. Sci., 40, pp 147-154.
- [6]- Douali L., Villemin D., Cherqaoui D.; (2004), "Exploring QSAR of non-nucleoside reverse transcriptase inhibitors by neural networks: TIBO derivatives", Int. Mol. Sci., 5, pp 48-55.
- [7] -Arab Chamjangali M., (2009), "Modeling of cytotoxicity data (CC<sub>50</sub>) of anti-HIV 1-[5-chlorophenyl] sulfonyl]-1H-pyrrole derivatives using calculated molecular descriptors and Levenberg- Marquardt artificial neural network", J.Chem. Bio. Drug. Des., 73, pp 456-465.
- [8]-Patel A., Gupta L., Karthikeyan C., Trivedi P.(2010), "QSAR studies on dihydro-alkoxy-benzyl-oxopyrimidines(DABOs) derivatives, as new series of potent, broad-spectrum non-nucleoside reverse transcriptase inhibitors" J. Current Pharmaceutical Research , 1, pp 19-25
- [9]- Chamjangali M. A., & Ashrafi M., (2013), " QSAR study of necroptosis inhibitory activities (EC<sub>50</sub>) of [1, 2, 3] thiadiazole and thiophene derivatives using Bayesian regularized artificial neural network and calculated descriptors", Medicinal Chemistry Research, 22(1), 392-400.

[10] عجم ز، (۱۳۹۲)، پایان نامه کارشناسی ارشد، "پیش بینی فعالیت ضد ایدز مشتقات غیر نوکلوزیدی تیوکربومات به روش جنگل های تصادفی"، دانشکده شیمی، دانشگاه صنعتی شاهرود.

[11]- Arab Chamjangali M., (2009), "Modeling of cytotoxicity data (CC<sub>50</sub>) of anti-HIV 1-[5-chlorophenyl] sulfonyl-1H-pyrrole derivatives using calculated molecular descriptors and Levenberg- Marquardt artificial neural network", J.Chem. Bio. Drug. Des., 73, pp 456-465.

[12]- Jalali-Heravi M., Asadollahi-Baboli, (2008), "QSAR analysis of platelet-derived growth inhibitors using GA-ANN and shuffling cross validation", J. QSAR. Comb. Sci., 6, pp 750-757

[13]- <http://www.wordiq.com/definition/chemometrics>

[14] - عرب چم جنگلی م،، (۱۳۸۶) "پیش بینی فعالیت دارویی ضد ایدز (سیتوتوکسیتی) مشتقات ۵-فنیل-۱-فنیل آمینو-۱-H-ایمیدازول به وسیله شبکه عصبی مصنوعی"، دانشگاه صنعتی شاهرود، گزارش طرح پژوهشی.

[15]- Arab Chamjangali M., Beglari M., Bagherian G., (2007), "Prediction of cytotoxicity data (CC<sub>50</sub>) of anti-HIV 5-phenyl-1-phenylamino-1H-imidazole derivatives by artificial neural network trained with Levenberg-Marquardt algorithm", J. Mol. Graph. Model, 26, pp 360-367.

[16]- بهرام م، خیامیان ت، عبدالهی ح، فاطمی م.ح، قاسمی ج، کمپانی زارع م، ناصری ع، همتی نژاد ب، (۱۳۹۴)، "مقدمه ای بر کمومتریکس"، چاپ اول، نشر پژوهشی نوآوران شریف، ۲۳۶، ص ۱۷۴-۱۷۲.

[17]- Jurs P. C. (2005) "Assessing the reliability of a QSAR model's predictions", J.Mol Graph Model, 23(6), pp 503.

[18]- شعیری ف، (۱۳۸۵)، پایان نامه کارشناسی ارشد، "پیشگویی فعالیت پیرازولوکوئینولین ها با استفاده از شبکه عصبی موجک و کاربرد روش RAFA در اندازه گیری همزمان آمپی سیلین و آموکسی سیلین"، دانشگاه علوم پایه دامغان.

[19]- Todeschini R. and Consonni V. (2000), "Handbook of Molecular Descriptors", Jonh Wiley, New York, pp98.

[20]- Daniel W. W., & Wayne W. D. (1995). "Biostatistics: a foundation for John WILEY. analysis in the health sciences".

[21]- بهرام م، خیامیان ت، عبدالهی ح، فاطمی م.ح، قاسمی ج، کمپانی زارع م، ناصری ع، همتی نژاد ب، (۱۳۹۴)، "مقدمه ای بر کمومتریکس"، چاپ اول، نشر پژوهشی نوآوران شریف، ۲۳۶، ص ۱۷۶-۱۷۷.

[22]- محمدرضایی م، (۱۳۸۹)، پایان نامه کارشناسی ارشد، "مطالعه کمی ساختار- فعالیت برخی از ترکیبات سولفونانیلید به عنوان گروهی جدید از داروهای ضد سرطان و ضد HIV برخی از ترکیبات"، دانشکده شیمی، دانشگاه صنعتی شاهرود.

[23] - کشتکار ه، (۱۳۹۳)، پایان نامه کارشناسی ارشد، "مطالعه کمی ساختار- خاصیت (QSPR) دانسیته مایعات

- یونی به صورت تابعی از دما و فشار با استفاده از روش‌های غیر خطی"، دانشکده شیمی، دانشگاه صنعتی شاهرود.
- [24]- فرشادفرع، (۱۳۸۰)، "اصول و روش‌های پیشرفته آماری (تجزیه رگرسیون)"، چاپ دوم، انتشارات طاق بستان
- [25] – <http://www.iranhiv.com/treatment>.
- [26] -هاگان ت. ه، دیموت ه، بیل م، "طراحی شبکه‌های مصنوعی"، کیا م، انتشارات کیان رایانه سبز.
- [27] -منهاج م، (۱۳۸۷)، "مبانی شبکه‌های عصبی (هوش محاسباتی)"، جلد اول، چاپ پنجم، مرکز نشر دانشگاه صنعتی امیرکبیر، تهران.
- [28] -بن کوپین، داورپناه، میرزایی ر، (۱۳۸۵)، "هوش مصنوعی"، انتشارات سخن گستر، مشهد.
- [29] -منهاج م، سیفی پور ن، (۱۳۷۷)، "کاربرد هوش محاسباتی در کنترل"، مرکز نشر دانشگاه صنعتی امیرکبیر.
- [30] - اشرفی م، (۱۳۸۹)، پایان نامه کارشناسی ارشد، "مطالعه کمی ساختار- فعالیت مشتقات تیوکربامات‌ها به عنوان دسته جدیدی از بازدارنده‌های غیرنوکلئوزیدی HIV"، دانشکده شیمی، دانشگاه صنعتی شاهرود.
- [31] -فردوسی م.ع، (۱۳۸۹)، پایان نامه کارشناسی ارشد، "پیش‌بینی ثابت‌های هنری بعضی از ترکیبات آلی با استفاده از روش‌های خطی و غیرخطی QSPR"، دانشکده شیمی، دانشگاه صنعتی شاهرود.
- [32]- Breiman L, (2001), "Random forests", Mach Learn, 45, pp 5-32
- [33]- عمادی جندقی ف، (۱۳۹۲)، پایان نامه کارشناسی ارشد، "کاربرد روش جنگل‌های تصادفی به عنوان ابزاری برای انتخاب متغیر و پیش‌بینی اندیس بازدارندگی از ترکیبات آلی به عنوان آلاینده محیط زیست"، دانشکده شیمی، دانشگاه صنعتی شاهرود.
- [34] -Hastie T, Tibshirani R, Friedman J, (2009), "The Elements of Statistical Learning Data Mining Inferences, and Prediction", 2<sup>nd</sup> Ed, Springer, New York.
- [35]- Genuer R, Poggi J M, Tuleau-Malot Ch, (2010), "Variable selection using random forests", Pattern Recogn Lett, 31, pp 2225-2236.
- [36]- Seo, B. I., Uchil, V. R., Okello, M., Mishra, S., Ma, X. H., Nishonov, M., ... & Nair, V. (2011). Discovery of a potent HIV integrase inhibitor that leads to a prodrug with significant anti-HIV activity. ACS medicinal chemistry letters, 2(12), 877-881.
- [37]- HyperChem7.0 Toronto, Canada: HyperCube Inc, <http://www.hyper.com>
- [38]- SPSS for windows Statistical package for IBM PC, SPSS Inc, <http://www.spss.com>.
- [39]- MATLAB 2009, the Math Work, Inc., Natick, MA, USA.

- [40]- Zupan J., Gasteiger J., (1993), "Neural networks for chemists an introduction", VCH Publishers, Weinheim, Germany.
- [41]- Ewing T., Baber J. C., & Feher M. (2006), "Novel 2D fingerprints for ligand-based virtual screening", *Journal of chemical information and modeling*, 46(6), 2423-2431.
- [42]- P. Broto, G. Moreau, C. Vandycke, (1984), "Molecular structures: pp 66– Perception, autocorrelation descriptor and SAR studies", *Eur. J. Med. Chem.*, 19, 70.
- [43]- Todeschini R., Gromatica P., (1997), "3D- modeling and prediction by WHIM descriptors. Part5. Theory development and chemical meaning of WHIM descriptors", *J. Quant. Struct., Acta., Relat.*, 16, pp 113-119.
- [44]- Douali L., Villemin D., Cherqaoui D.; (2003), "Neural networks: Accurate nonlinear QSAR model for HEPT derivatives", *J. Chem. Inf. Comput. Sci.*, 43, pp 1200-1207.
- [45]- chen Z., Zhang Y., Fu W., (2010), "QSAR study of carboxylic acid derivatives as HIV-1 integrase inhibitors", *Eur. J. Med. Chem.*, 45, pp 3970-3980.
- [46]- Xu, Z., Ba, M., Zhou, H., Cao, Y., Tang, C., Yang, Y., ... & Zhu, L. (2014). 2, 4, 5-Trisubstituted thiazole derivatives: a novel and potent class of non-nucleoside inhibitors of wild type and mutant HIV-1 reverse transcriptase. *European journal of medicinal chemistry*, 85, 27-42.
- [47]- Crisan L., Iliescu S., & Funar-Timofei S.,(2016), "Structure-flammability relationship study of phosphoester dimers by MLR and PLS", *Polímeros, (AHEAD)*, 0-0.
- [48]- Spallarossa A., Cesarini S., Ranise A., et al.,(2009), "Parallel synthesis," molecular modelling and further structure–activity relationship studies of new acylthiocarbamates as potent non-nucleoside HIV-1 reverse transcriptase inhibitors". *European Journal of Medicinal Chemistry*, 44(5): p. 2190-2201.

## Abstract

In the first section of this study, quantitative structure – activity relationship (QSAR) models were constructed for predicting the anti-HIV activity of some pyridinone derivatives. The group contribution and the stepwise multiple linear regression methods were used for variable selection. Then the best descriptors were applied as input for QSAR modeling using random forest (RF) and artificial neural network (ANN) techniques. In modelling methods based on random forest the data set was divided into training set and test set including 43 and 10 compounds, respectively. Effective parameters of RF method including  $N_{tree}$  and  $Mtry$  were optimized using the training set. For ANN modeling the data set was divided into training set, validation set and test set including 37, 8 and 8 compounds, respectively. ANN model was optimized too and the best model was selected. The validation study of the RF and ANN models was performed by prediction of the activity of the test set compounds and the correlation coefficients for RF and ANN model are 0.929 and 0.935 respectively.

In the second section, RF and ANN models were constructed and the ability of prediction of anti-HIV activities for some of thiazole derivatives were investigated. The data set was divided into training and test set containing 18 and 6 chemical, respectively. The stepwise regression method was used to select the most important descriptors. The validation study of the RF and ANN models was performed using test set compounds. The correlation coefficients obtained for test set were 0.982 and 0.896 for RF and ANN respectively. The results obtained showed proper prediction power of the proposed models.

**Keyword:** QSAR, Anti-HIV, Artificial neural network, Random forest, group contribution



**Shahrood University of Technology**

**Faculty of Chemistry**

**MSc thesis in Analytical Chemistry**

**Quantitative structure-activity relationship study  
of the activity of some pharmaceutical compounds  
using Random forest**

**By: Atiyeh Asali**

Supervisors:

**Dr. Nasser Goudarzi**

**Dr. Mansour Arab Chamjangali**

September 2016