





دانشگاه صنعتی شاهرود

دانشکده شیمی

رساله دکتری شیمی تجزیه

به کارگیری تکنیک‌های رگرسیونی جریمه‌شده به عنوان روش‌های جدید انتخاب متغیر

در مطالعات روابط کمی ساختار-فعالیت (QSAR) و ساختار-خاصیت (QSPR)

نگارنده:

زینب مظفری

استاد راهنما:

دکتر منصور عرب‌چم‌جنگلی

اساتید مشاور:

دکتر محمد آرشی

دکتر ناصر گودرزی

به پاس تعبیر عظیم و انسانی‌شان از کلمه ایشار، به پاس عاطفه سرشار و گرمای  
امیدبخش وجودشان که در این سردترین روزگار ان بهترین پشتیبان است، به پاس  
قلب‌های بزرگشان که فریاد رس است و به پاس محبت‌های بی‌دیریشان که هرگز فروکش  
نمی‌کند،

این رساله را تقدیم می‌کنم:

به استوارترین تکیه‌گاهم، دستان پر مهر پدرم

به همراه‌ترین نگاه زندگیم، چشمان پر امید مادرم

به بهترین هم‌بازی دوران کودکی و بهترین حامی دوران بزرگسالی، برادر

عزیزم، ابراهیم

به انگیزه، امید و چراغ همیشه تابان زندگیم، خواهرمانزینم، فروزان

به نو رود باغ زندگیان، مرغان

و به تمام کسانی که نیروی حرکت به سوی مقصود را در وجود من زنده کردند

هرچه آموختم در مکتب عشق شما آموختم و هرچه بگوختم قطره‌ای از دریای بی

کران مهربانی‌تان را سپاس تو انم بگویم.

امروز هستی‌ام به امید شماست و فردا کلید باغ بهشتم رضای شما

ره آوردی کران سنگ تراز این رساله نداشتم تا به خاک پایتان شار کنم، باشد که

حاصل تلاشم نسیم کوزه غبار هستی‌تان را برزاید.

بوسه بردستان پر مهرتان، خدا پشت و پناهتان



## قدر دانی و تشکر:

حمد و ثنای خدای متعال را که به فضلش آدمی را اندیشه و عمل عطا فرمود تا بجاود و بداند و از علم و آگاهی برخوردار شود و حتی خویش را روشنی بخشد. پس از ارادت خاضعانه به درگاه خداوند بی همتا لازم است از انسانهای عاشق همواره در تفکر و شناخت، آشنایی که زندگی خویش را وقف علم و شریعت کردند قدر دانی نمایم. در اینجا لازم میدانم که از زحمات و راهنماییهای خردمندان و بی دریغ استاد راهنمای گرانقدرم، جناب آقای دکتر منصور عرب تقدیر و تشکر نمایم، که موفقیتهای پژوهشی و آموزشی بنده حاصل زحمات ایشان در باروری و بروز رسانی اطلاعات اینجانب می باشد. در تدوین این رمانورد خود را مدیون راهنمایی و رهنمودهای ارزنده و مساعدت صمیمانه استاد مشاور ارجمندم، جناب آقای دکتر محمد آرش می دانم، که نقش بسزایی را در به پایان رساندن و موفقیت این اثر داشته اند. به راستی که انجام این رساله بدون نظرهای صائب، پیکیریهای دلسوزانه و تشویقهای امیدبخش این دو عزیز میسر نبود. تشکر و قدر دانی خود را از جناب آقای دکتر ناصر کوردزی استاد و مشاور ارجمند اعلام می دارم که الطاف بی شائبه ایشان در طول مدت تحصیل شامل حال بنده شد. از استاد بزرگوارم جناب آقای دکتر قدوسی باقریان، برای تمام زحمات بی دریغشان و زحمت دایمی رساله، ممنونم. از اساتید بزرگوار، جناب آقای دکتر فاطمی و جناب آقای دکتر روزبه که زحمت دایمی و تصحیح این رساله را بر عهده داشتند کمال سپاس را دارم. از خداوند بزرگ توفیق، طول عمر، سلامتی و موفقیت را برای همه اساتید بزرگوار خواهانم.

از بهکلاسی، هم اتاقی، هم دانشگاهی، هم آزمایشگاهی، دوست و خواهر عزیزم خانم بهاره عربستانی که در طی این سال ها همراه همیشگی روزهای سخت و خاطر هساز ایام خوش در شهر غربت بود، صمیمانه تقدیر می کنم و برای ایشان از خداوند بزرگ عاقبت به خیری و موفقیت آرزو مندم.

از همزمان همیشگی ام خانم غزاله نجفی عرب، خانم یاسین قبری، خانم نسیرین مماندوست، خانم آرزو افروز، خانم دکتر اشرفی، خانم دکتر دوستی، خانم محدثه لطنی، آقای مهندس داوود نادعلی، آقای مهندس مومنی و آقای مهندس یزدانی که اوقات خوشی را در کنار هم سپری کرده ایم و در این مدت، همیشه بنده را مورد لطف و مرحمت خود قرار داده اند، کمال تقدیر و تشکر را دارم. در نهایت بر خود لازم می دانم تا از دوستان عزیزم در آزمایشگاه شیمی تجزیه، شیمی آلی و شیمی معدنی، اعضای هیات علمی دانشکده شیمی به خصوص آقای دکتر بهرامیان، آقای دکتر میرزایی، سرکار خانم دکتر مصدرا لامور و سرکار خانم دکتر کلاتر، کارشناسان و کارکنان دانشکده شیمی، و کارشناسان تحصیلات تکمیلی به خصوص آقای مهندس خانعلی زاده، که در دوره ۷ ساله تحصیلی کارشناسی ارشد و دکتری، همواره اینجانب را راهنمایی و همراهی نموده اند، سپاسگزارم.

## تعهد نامه

اینجانب زینب مظفری دانشجوی دوره دکتری شیمی تجزیه دانشکده شیمی دانشگاه صنعتی شاهرود، نویسنده رساله به کارگیری تکنیک‌های رگرسیون جرمی به‌عنوان روش‌های جدید انتخاب متغیر در مطالعات روابط کمی ساختار-فعالیت (QSAR) و ساختار-خاصیت (QSPR) تحت راهنمایی آقای دکتر منصور عرب چم جنگلی و مشاوره آقای دکتر محمد آرشی و آقای دکتر ناصر گودرزی متعهد می‌شوم:

- 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 می‌باشد 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 بوده‌اند در مقالات مستخرج از پایان نامه رعایت می‌گردد.
  - 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.
- تاریخ امضای دانشجو**

مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج، کتاب، برنامه‌های رایانه‌ای، نرم افزارها و تجهیزات ساخته شده است) متعلق به دانشگاه صنعتی شاهرود می‌باشد. 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 نمی‌باشد.

## چکیده

اهداف مهم این رساله، استفاده از روش‌های انتخاب متغیر جریمه‌شده برای انتخاب مؤثرترین توصیف‌کننده‌ها و جفت کردن روش‌های جریمه‌شده با روش مدل‌سازی غیر خطی شبکه عصبی مصنوعی (ANN) است. در این راستا، از روش‌های مختلفی برای ایجاد ارتباط بین توصیف‌کننده‌های منتخب و پاسخ هدف استفاده شد. بنابراین، در مطالعه اول از ترکیب انحراف قدر مطلق به‌طور هموار بریده شده (SCAD) و ANN به عنوان یک رویکرد جدید (SCAD-LM-ANN) در مطالعات کمی ساختار-فعالیت (QSAR) استفاده شد. SCAD-LM-ANN از مزایای ذاتی مفید روش انقباضی SCAD در کاهش داده‌هایی با ابعاد بالا قبل از روش مدل‌سازی استفاده می‌کند. عملکرد مدل SCAD-LM-ANN با ایجاد ارتباط بین توصیف‌کننده‌های به‌دست‌آمده از دراگون و فعالیت‌های دارویی برای مجموعه‌ای از مشتقات تیواستامید/استانیلید به عنوان مهارکننده‌های HIV مورد بررسی قرار گرفت. روش SCAD با اجرای ارزیابی تقاطعی  $10$ -fold بر روی مجموعه داده‌ها در غیاب ترکیبات مجموعه آزمون، اجرا شد.  $11$  توصیف‌کننده در  $\lambda$  با کم‌ترین خطای ارزیابی تقاطعی ( $\lambda_{\min}$ ) انتخاب شدند. توصیف‌کننده‌های منتخب به‌عنوان ورودی ANN مورد استفاده قرار گرفتند. تمام پارامترهای مؤثر بر عملکرد مدل، بهینه شدند و مدل SCAD-LM-ANN با معماری  $5-5-1$  به‌عنوان مدل QSAR بهینه انتخاب شد. چندین پارامتر آماری برای ارزیابی مدل در نظر گرفته شد. نتایج حاصل، تعمیم‌پذیری و قدرت پیش‌بینی مدل SCAD-LM-ANN پیشنهادی را اثبات می‌کند. با توجه به رابطه ایجاد شده در مدل QSAR برتر، مشتقات جدیدی طراحی و به عنوان مهارکننده‌های فعال HIV برای مطالعات پیش‌تر پیشنهاد شدند. صحت ترکیبات پیشنهادی با تجزیه و تحلیل برهم‌کنش‌های گیرنده-لیگاند (LR) به‌دست آمده از مطالعات داکینگ مولکولی، مورد مطالعه و تأیید قرار گرفت. در بخش دوم این رساله، رویکرد جدیدی به عنوان ترکیبی از حداقل قدر مطلق انقباض و عملگر انتخاب‌کننده سازگار (ALASSO) و ANN برای ساخت مدل مهارکننده‌های پروتئاز ۳-شبه کیموتریپسین ( $3CL^{pro}$ ) به‌عنوان ترکیبات قوی SARS CoV-2 معرفی شد. با توجه به اهمیت این بیماری از سال ۲۰۱۹، توصیه و طراحی ترکیبات فعال جدید بسیار مهم است. پس از ارزیابی صحت و اعتبار مدل توسعه‌یافته ALASSO-LM-ANN، ترکیبات جدیدی با استفاده از توصیف‌کننده‌های مؤثر پیشنهاد شدند و فعالیت دارویی ترکیبات جدید پیش‌بینی شد. بررسی برهم‌کنش‌های LR نیز با استفاده از مطالعه داکینگ مولکولی انجام شد. خواص PK و قانون

پنج لیپینسکی برای تمام ترکیبات پیشنهادی محاسبه شد و ترکیبات جدید پیشنهادی دارای خواص دارویی قابل قبولی هستند. در مطالعه سوم، ترکیب حداقل انحراف مطلق - حداقل قدر مطلق انقباض و عملگر انتخاب کننده (LAD-LASSO) به عنوان یک روش انتخاب متغیر جدید برای مطالعات QSAR مبتنی بر ANN معرفی شد. مدل ANN همراه با روش انتخاب متغیر کارآمد LAD-LASSO برای پیش بینی فعالیت دارویی سه مجموعه داده از ترکیبات شیمیایی، مورد ارزیابی قرار گرفت. از روش انتخاب متغیر LAD-LASSO استفاده شد و توصیف کننده‌هایی که بیشترین ارتباط را با فعالیت‌های دارویی داشتند انتخاب شدند. توصیف کننده‌های منتخب به عنوان ورودی ANN تعریف شدند و مدل‌های طراحی شده، بهینه شدند. فعالیت دارویی ترکیبات مجموعه آموزش با استفاده از مدل‌های بهینه ANN پیش بینی شد. ضرایب تعیین  $(R^2)$  برای داده‌های آزمون در سه مجموعه داده برابر با ۰/۸۷، ۰/۸۴ و ۰/۸۷ بود. دامنه کاربرد و آزمون  $-Y$  تصادفی نیز کارایی مدل‌های توسعه یافته را ثابت کردند. در نهایت، مدل‌های QSAR ایجاد شده برای پیشنهاد ترکیبات شیمیایی جدید و قوی با اصلاح ساختاری مولکول‌های ضعیف در هر سه مجموعه داده مورد استفاده قرار گرفتند. مقادیر پاسخ ترکیبات پیشنهادی جدید با استفاده از مدل‌های ANN بهینه پیش بینی شدند. بر اساس اطلاعات LR، وجود برهم کنش‌های مختلف آب‌دوست و آب‌گریز در جایگاه فعال گیرنده نشان‌دهنده پتانسیل بالای ترکیبات شیمیایی در برقراری اتصال پایدار است. در آخرین کار این مطالعه، ترکیبی از SCAD و ANN در رابطه کمی ویژگی ساختاری - شاخص‌های بازداری (QSRR) (RIs) استفاده شد. روش SCAD پیشنهادی قبل از استفاده از روش قدرتمند مدل‌سازی ANN، ابعاد داده‌ها را کاهش می‌دهد. کارایی روش‌های SCAD-ANN با ساخت یک مدل QSRR بین مؤثرترین توصیف کننده‌های مولکولی و RI برای دو مجموعه از ترکیبات آلی فرار (VOCs) ارزیابی شد. روش SCAD برای داده‌های آموزشی اعمال شد و توصیف کننده‌های مؤثر در  $\lambda_{min}$  انتخاب شدند و به عنوان ورودی‌های روش مدل‌سازی ANN تعریف شدند. تمام پارامترهای ANN به‌طور هم‌زمان بهینه شدند. نتایج به‌دست آمده نشان می‌دهد که مدل‌های QSRR ساخته شده از قدرت پیش بینی قابل قبولی برخوردار است.

**کلمات کلیدی:** QSAR، QSPR، سرطان، HIV، Coronavirus، SCAD، ALASSO، LAD-LASSO.

ANN، اتصال مولکولی

## مقالات چاپ شده در مجلات خارجی JCR

- 1- **Mozafari, Z.**, Arab Chamjangali, M., Arashi, M., & Goudarzi, N. (2021). Performance of smoothly clipped absolute deviation as a variable selection method in the artificial neural network- based QSAR studies. *Journal of Chemometrics*, e3338.
- 2- **Mozafari, Z.**, Chamjangali, M. A., Arashi, M., & Goudarzi, N. (2021). Suggestion of active 3-chymotrypsin like protease (3CL<sup>Pro</sup>) inhibitors as potential anti-SARS-CoV-2 agents using predictive QSAR model based on the combination of ALASSO with an ANN model. *SAR and QSAR in Environmental Research*, 32(11), 863-888.
- 3- **Mozafari, Z.**, Chamjangali, M. A., Arashi, M., & Goudarzi, N. (2021). QSRR models for predicting the retention indices of VOCs in different datasets using an efficient variable selection method coupled with artificial neural network modeling: ANN-based QSPR modeling. *journal of iranian chemical society*. Accepted manuscript on 18 Dec. 2021.
- 4- **Mozafari, Z.**, Chamjangali, M. A., Arashi, M., & Goudarzi, N. (2022). Application of the LAD-LASSO as a dimensional reduction technique in the ANN-based QSAR study: Discovery of potent inhibitors using molecular docking simulation. *Chemometrics and Intelligent Laboratory Systems*, 104510.

## مقالات ارائه شده در کنفرانس‌های معتبر داخلی به شکل سخنرانی و پوستر

5. **Mozafari, Z.**, Arab Chamjangali, M., Arashi, M., & Goudarzi, N. Hybrid QSPR models for the prediction of the linear retention index of volatile compounds in flour, 7<sup>th</sup> Iranian Biennial Chemometrics Seminar, 30-31 Oct. 2019, Shahrood University of Technology, Shahrood, Iran, **(Poster)**
6. **Mozafari, Z.**, Arab Chamjangali, M., Arashi, M., & Goudarzi, N. Application of a new hybrid of SCAD - artificial neural network in QSAR study of HIV inhibitors, 7<sup>th</sup> Iranian Biennial Chemometrics Seminar, 30-31 Oct. 2019, Shahrood University of Technology, Shahrood, Iran, **(Oral)**
7. **Mozafari, Z.**, Arab Chamjangali, M., Arashi, M., & Goudarzi, N. Application of the LAD-LASSO as a dimensional reduction technique in the ANN-based QSAR study: Discovery of potent inhibitors using molecular docking simulation, 10<sup>th</sup> Theoretical and Computational Chemistry Workshop, 23 – 25 Nov 2021, Sharif University of Technology, Tehran, Iran, **(Oral)**

## فهرست مطالب

فصل اول: مقدمه	۱
۱-۱ مقدمه	۲
۱-۲ کمومتریکس	۳
۱-۳ رابطه کمی ساختار- فعالیت (QSAR)	۴
۱-۴ رابطه کمی ساختار- خاصیت (QSPR)	۶
۱-۵ اهمیت ساخت مدل‌های ارتباط کمی ساختار- فعالیت/ ویژگی (QSAR/QSPR)	۷
۱-۵-۱ مراحل ساخت مدل‌های QSAR/QSPR	۸
۱-۵-۲ جمع‌آوری و انتخاب مجموعه داده‌ها	۱۰
۱-۵-۳ رسم و بهینه‌سازی ساختارهای شیمیایی	۱۰
۱-۵-۴ استخراج توصیف‌کننده‌ها	۱۱
۱-۵-۵ پیش‌پردازش توصیف‌کننده‌های محاسبه شده	۱۱
۱-۵-۶ انتخاب توصیف‌کننده‌های مؤثر	۱۲
۱-۵-۶-۱ انتخاب متغیر به روش رگرسیون خطی چندگانه	۱۳
۱-۵-۶-۲ LASSO	۱۴
۱-۵-۶-۳ SCAD	۱۶
۱-۵-۶-۴ ALASSO	۱۷
۱-۵-۶-۵ LAD-LASSO	۱۸
۱-۵-۶-۶ ارزیابی توصیف‌کننده‌های منتخب	۱۹
۱-۵-۷ مدل‌سازی	۲۰
۱-۵-۷-۱ شبکه عصبی مصنوعی	۲۰

- ۱- ۵- ۷- ۲ چیدمان توصیف‌کننده‌های منتخب به‌عنوان ورودی مدل شبکه عصبی مصنوعی..... ۲۳
- ۱- ۵- ۸ ارزیابی مدل..... ۲۵
- ۱- ۵- ۸- ۱ بررسی مدل با استفاده از نتایج پیش‌بینی شده برای مجموعه آزمون ..... ۲۵
- ۱- ۵- ۸- ۲ ارزیابی مدل با استفاده از تکنیک رد مرحله‌ای تک تک ..... ۲۶
- ۱- ۵- ۸- ۳ نمودار باقی‌مانده‌های استاندارد شده ..... ۲۶
- ۱- ۵- ۸- ۴ پارامترهای آماری ..... ۲۷
- ۱- ۵- ۸- ۵ دامنه کاربرد مدل ..... ۳۰
- ۱- ۵- ۸- ۶ آزمون Y- تصادفی ..... ۳۱
- ۱- ۵- ۸- ۷ بررسی مشارکت توصیف‌کننده‌های منتخب در شبکه عصبی ..... ۳۱
- ۱- ۶ شبیه‌سازی داکینگ مولکولی ..... ۳۲
- ۱- ۷- ۱ مراحل اجرای داکینگ مولکولی ..... ۳۳
- ۱- ۷- ۱ آماده‌سازی پروتئین ..... ۳۴
- ۱- ۷- ۲ ساختن لیگاند ..... ۳۵
- ۱- ۷- ۳ تنظیم کردن جعبه شبکه‌ای ..... ۳۵
- ۱- ۷- ۴ گزینه‌های داکینگ ..... ۳۶
- ۱- ۷- ۵ انجام محاسبه داکینگ ..... ۳۶
- ۱- ۷- ۶ آنالیز و تحلیل نتایج ..... ۳۶
- ۱- ۷- ۷ کاربردهای داکینگ مولکولی ..... ۳۷
- ۱- ۸- ۱ اهمیت و ضرورت مدل‌سازی QS(A/P)R ترکیبات شیمیایی با استفاده از مدل‌های شبکه عصبی توسعه یافته با توصیف‌کننده‌های منتخب روش‌های جریمه‌ای و مروری بر کارهای انجام شده ..... ۳۸
- ۱- ۸- ۱ اهمیت پیش‌بینی فعالیت دارویی بازدارنده‌های ایدز با استفاده از مدل SCAD-ANN ..... ۴۰

۴۴.....	۸-۱-۲ اهمیت پیش‌بینی فعالیت دارویی بازدارنده‌های SARS-COV-2 با استفاده از مدل ALASSO-ANN
۴۵.....	۸-۱-۳ اهمیت پیش‌بینی فعالیت دارویی برخی از بازدارنده‌های ایدز و سرطان با استفاده از مدل LAD-LASSO-ANN
۴۸.....	۸-۱-۴ اهمیت پیش‌بینی شاخص بازدارندگی ترکیبات آلی فرار با استفاده از مدل SCAD-ANN
۵۰.....	۹-۱-۹ مروری بر تحقیقات انجام شده در مورد به‌کارگیری روش‌های انتخاب متغیر جریمه‌ای در ساخت مدل‌های QSAR/QSPR
۵۴.....	۱۰-۱ نوآوری تحقیق.....
۵۵.....	<b>فصل دوم: بخش تجربی</b>
۵۶.....	۲-۱ معرفی نرم‌افزارهای مورد استفاده برای مدل‌سازی QSAR
۵۶.....	۲-۱-۱ نرم‌افزار هایپرکم
۵۶.....	۲-۱-۲ نرم‌افزار دراگون
۵۷.....	۲-۱-۳ نرم‌افزار SPSS
۵۷.....	۲-۱-۴ نرم‌افزارهای R و R-studio
۵۸.....	۲-۱-۵ نرم‌افزار متلب
۵۹.....	۲-۱-۶ نرم‌افزار Origin Lab
۵۹.....	۲-۱-۷ نرم‌افزار اتوداک
۵۹.....	۲-۱-۸ نرم‌افزار ویورلایت
۶۰.....	۲-۱-۹ نرم‌افزار بیوویا دی اس
۶۰.....	۲-۱-۱۰ نرم‌افزار VMD



۲-۲	پیش‌بینی فعالیت دارویی برخی از مشتقات استانیلید/استامید به‌عنوان بازدارنده‌های ایدز با استفاده از مدل SCAD-ANN	۶۱
۲-۲-۱	مقدمه	۶۱
۲-۲-۲	مجموعه داده‌ها	۶۲
۲-۲-۳	رسم و بهینه‌سازی ساختار Arylazolythioacetamide/acetanilide ها	۶۵
۲-۲-۴	استخراج توصیف‌کننده‌ها	۶۵
۲-۲-۵	پیش‌پردازش و انتخاب توصیف‌کننده‌های مؤثر	۶۵
۲-۲-۶	مدل‌سازی شبکه عصبی با استفاده از توصیف‌کننده‌های منتخب SCAD	۶۸
۲-۲-۷	ارزیابی مدل SCAD-LM-ANN	۷۱
۲-۲-۷-۱	ارزیابی مدل SCAD-LM-ANN با استفاده از پیش‌بینی داده‌های مجموعه آزمون	۷۱
۲-۲-۷-۲	ارزیابی مدل SCAD-ANN با استفاده از روش رد مرحله‌ای تک تک	۷۳
۲-۲-۷-۳	ارزیابی مدل SCAD-LM-ANN با استفاده از پارامترهای آماری	۷۷
۲-۲-۷-۴	ارزیابی مدل SCAD-LM-ANN با استفاده از دامنه کاربرد	۷۸
۲-۲-۷-۵	ارزیابی مدل SCAD-LM-ANN با استفاده از آزمون Y-تصادفی	۸۰
۳-۲	پیش‌بینی فعالیت دارویی برخی از مشتقات ۳-(3CL <sup>Pro</sup> ) chymotrypsin like protease به‌عنوان بازدارنده‌های SARS-COV-2 با استفاده از مدل ALASSO-ANN	۸۲
۳-۲-۱	مقدمه	۸۲
۳-۲-۲	مجموعه داده‌ها	۸۵
۳-۲-۳	رسم و بهینه‌سازی ساختار بازدارنده‌های 3CL <sup>Pro</sup>	۸۹
۳-۲-۴	استخراج توصیف‌کننده‌ها	۸۹
۳-۲-۵	پیش‌پردازش و انتخاب توصیف‌کننده‌های مؤثر	۸۹

- ۹۲-۳-۲ مدل سازی شبکه عصبی با استفاده از توصیف کننده های منتخب ALASSO.....
- ۹۴-۳-۲ ارزیابی مدل ALASSO-LM-ANN.....
- ۹۴-۳-۲-۱ ارزیابی مدل ALASSO-LM-ANN با استفاده از مجموعه آزمون.....
- ۹۴-۳-۲-۲ ارزیابی مدل ALASSO-LM-ANN با پیش بینی  $pIC_{50}$  تمام ترکیبات مجموعه داده با استفاده از روش رد مرحله ای تک تک.....
- ۹۶-۳-۲-۳ ارزیابی مدل ALASSO-LM-ANN با استفاده از پارامترهای آماری.....
- ۱۰۱-۳-۲-۴ ارزیابی ALASSO-LM-ANN با استفاده از دامنه کاربرد.....
- ۱۰۳-۳-۲-۵ ارزیابی مدل ALASSO-LM-ANN با استفاده از آزمون Y-تصادفی.....
- ۱۰۵-۲-۴ پیش بینی فعالیت دارویی برخی از بازدارنده های ایدز و سرطان با استفاده از مدل LAD-LASSO-ANN.....
- ۱۰۵-۴-۲-۱ مقدمه.....
- ۱۰۷-۴-۲-۲ مجموعه داده ها.....
- ۱۱۳-۴-۲-۳ رسم و بهینه سازی ساختار ترکیبات شیمیایی مجموعه داده های متفاوت.....
- ۱۱۳-۴-۲-۴ استخراج توصیف کننده های ساختاری.....
- ۱۱۴-۴-۲-۵ پیش پردازش و انتخاب توصیف کننده های مؤثر.....
- ۱۲۰-۴-۲-۶ مدل سازی شبکه عصبی با استفاده از توصیف کننده های منتخب LAD-LASSO.....
- ۱۲۲-۴-۲-۷ ارزیابی مدل LAD-LASSO-LM-ANN هر سه مجموعه داده.....
- ۱۲۳-۴-۲-۱ ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از پیش بینی داده های مجموعه آزمون.....
- ۱۲۳-۴-۲-۲ ارزیابی مدل LAD-LASSO-LM-ANN با پیش بینی فعالیت دارویی تمام ترکیبات مجموعه داده ها با استفاده از روش رد مرحله ای تک تک.....
- ۱۲۷-۴-۲-۳ ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از پارامترهای آماری.....
- ۱۳۴-۴-۲-۴ ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از دامنه کاربرد.....
- ۱۳۶-۴-۲-۴-۷ ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از دامنه کاربرد.....
- ۱۳۸-۴-۲-۵ ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از آزمون Y-تصادفی.....

۱۴۱	۲-۵ پیش‌بینی شاخص بازداری برخی از ترکیبات آلی فرار با استفاده از مدل SCAD-ANN
۱۴۱	۲-۵-۱ مقدمه
۱۴۳	۲-۵-۲ مجموعه داده‌ها
۱۴۹	۲-۵-۳ رسم و بهینه‌سازی ساختار ترکیبات آلی فرار مجموعه داده‌های متفاوت
۱۴۹	۲-۵-۴ استخراج توصیف‌کننده‌های ساختاری
۱۴۹	۲-۵-۵ پیش‌پردازش و انتخاب توصیف‌کننده‌های مؤثر
۱۵۴	۲-۵-۶ مدل‌سازی شبکه عصبی با استفاده از توصیف‌کننده‌های منتخب SCAD
۱۵۶	۲-۵-۷ ارزیابی مدل SCAD-ANN
۱۵۷	۲-۵-۷-۱ ارزیابی مدل SCAD-ANN با استفاده از پیش‌بینی مجموعه داده‌های آزمون
	۲-۵-۷-۲ ارزیابی مدل SCAD-ANN با پیش‌بینی شاخص‌های بازداری تمام ترکیبات مجموعه داده‌های A و B با استفاده از روش رد مرحله‌ای تک تک
۱۶۰	
۱۶۷	۲-۵-۷-۳ ارزیابی مدل SCAD-ANN با استفاده از پارامترهای آماری
۱۶۹	۲-۵-۷-۴ ارزیابی مدل SCAD-ANN با استفاده از دامنه کاربرد
۱۷۱	۲-۵-۷-۵ ارزیابی مدل SCAD-ANN با استفاده از آزمون Y-تصادفی
۱۷۳	<b>فصل سوم: نتیجه‌گیری و آینده‌نگری</b>
۱۷۴	۳-۱ نتیجه‌گیری نهایی مدل‌های توسعه یافته QSAR/QSPR
۱۷۵	۳-۱-۱ تجزیه و تحلیل توصیف‌کننده‌های مدل SCAD-LM-ANN برای مجموعه بازدارنده‌های ایدز
۱۷۵	۳-۱-۱-۱ محاسبه سهم مشارکت هر توصیف‌کننده در مدل SCAD-ANN
	۳-۱-۱-۲ بررسی رابطه بین توصیف‌کننده‌های استفاده شده در مدل نهایی (SCAD-LM-ANN) و فعالیت دارویی
۱۷۶	ترکیبات مورد مطالعه
۱۷۸	۳-۱-۱-۳ پیشنهاد ترکیبات جدید با فعالیت دارویی مناسب با استفاده از مدل SCAD-LM-ANN ارائه شده
۱۷۹	۳-۱-۱-۴ مطالعه داکینگ مولکولی

- ۳-۱-۲ تجزیه و تحلیل توصیف‌کننده‌های مدل ALASSO-ANN برای بازدارنده‌های SARS-COV-2..... ۱۸۸
- ۳-۱-۲-۱ محاسبه سهم مشارکت هر توصیف‌کننده در مدل ALASSO-ANN..... ۱۸۸
- ۳-۱-۲-۲ بررسی رابطه بین توصیف‌کننده‌های استفاده شده در مدل نهایی (ALASSO-LM-ANN) و فعالیت دارویی ترکیبات مورد مطالعه..... ۱۸۹
- ۳-۱-۲-۳ کاربرد مدل ALASSO-LM-ANN در طراحی و پیشنهاد ترکیبات فعال با اثر ضد کووید-۱۹..... ۱۹۲
- ۳-۱-۲-۴ مطالعه داکینگ مولکولی بازدارنده‌های 3CL<sup>pro</sup>..... ۱۹۳
- ۳-۱-۳ تجزیه و تحلیل توصیف‌کننده‌های مدل LAD-LASSO-LM-ANN برای مجموعه داده‌های ضد ایدز و ضد سرطان..... ۲۱۰
- ۳-۱-۳-۱ محاسبه سهم مشارکت هر توصیف‌کننده در مدل LAD-LASSO-LM-ANN..... ۲۱۰
- ۳-۱-۳-۲ بررسی میزان و چگونگی تأثیر توصیف‌کننده‌ها بر فعالیت دارویی ترکیبات مجموعه داده‌های ضد ایدز و کاربرد مدل LAD-LASSO-LM-ANN ارائه شده در پیشنهاد ترکیبات جدید..... ۲۱۲
- ۳-۱-۳-۳ مطالعه داکینگ مولکولی بازدارنده‌های ایدز..... ۲۱۵
- ۳-۱-۳-۴ بررسی میزان و چگونگی تأثیر توصیف‌کننده‌ها بر فعالیت دارویی ترکیبات مجموعه داده‌های ضد سرطان کارسینوم کولورکتال و ریه و کاربرد مدل LAD-LASSO-LM-ANN ارائه شده در پیشنهاد ترکیبات جدید..... ۲۲۶
- ۳-۱-۴ تحلیل توصیف‌کننده‌های مدل SCAD-ANN برای ترکیبات آلی فرار..... ۲۳۸
- ۳-۱-۴-۱ محاسبه سهم مشارکت هر توصیف‌کننده در مدل SCAD-ANN..... ۲۳۸
- ۳-۱-۴-۲ بررسی رابطه بین توصیف‌کننده‌های استفاده شده در مدل SCAD-ANN و شاخص بازدارندگی ترکیبات مورد مطالعه..... ۲۳۹
- ۳-۲ نتیجه‌گیری نهایی..... ۲۴۳
- ۳-۳ آینده‌نگری..... ۲۵۱
- ۳-۴ منابع..... ۲۵۲

## فهرست شکل‌ها:

- شکل ۱-۱ مراحل ساخت، توسعه و ارزیابی مدل‌های QSAR/QSPR ..... ۹
- شکل ۲-۱ ساختار کلی HIV [۵۴] ..... ۴۱
- شکل ۳-۱ جدول زمانی کوتاهی از توسعه داروهای HIV، داروهای شناسایی شده برای RT با خط چین (NRTIs) و خط (NNRTIs) جدا شده‌اند [۵۴] ..... ۴۳
- شکل ۴-۱ شکل بازدارنده‌های بالقوه PI3K. A ورتمانین، B: LY294002 و C کوئرستین ..... ۴۷
- شکل ۱-۲ ساختار پایه ترکیبات Arylazolythioacetamide/acetanilide مورد مطالعه ..... ۶۲
- شکل ۲-۲ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب SCAD ..... ۶۷
- شکل ۳-۲ نمودار مقادیر VIF توصیف‌کننده‌های منتخب SCAD ..... ۶۷
- شکل ۴-۲ نمودار تغییرات مقادیر پیش‌بینی شده  $pEC_{50}$  به‌وسیله مدل SCAD-LM-ANN در شرایط بهینه در مقابل مقادیر تجربی برای داده‌های مجموعه آزمون ..... ۷۳
- شکل ۵-۲ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌ها بر اساس تکنیک LOO در مقابل مقادیر تجربی ..... ۷۶
- شکل ۶-۲ نمودار باقی‌مانده‌های حاصل از پیش‌بینی فعالیت دارویی ترکیبات با استفاده از تکنیک LOO و مقادیر تجربی برحسب مقادیر تجربی ..... ۷۶
- شکل ۷-۲ دامنه کاربرد مدل SCAD-LM-ANN، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $\sigma \pm 3$  و  $h^*$  است ..... ۸۰
- شکل ۸-۲ نمودار مقادیر  $R^2$  به‌دست آمده در آزمون Y-تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای Y-تصادفی و پیش‌بینی فعالیت ترکیبات مجموعه آزمون به‌وسیله مدل SCAD-LM-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه ..... ۸۱
- شکل ۹-۲ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب روش ALASSO ..... ۹۱
- شکل ۱۰-۲ نمودار مقادیر VIF توصیف‌کننده‌های منتخب روش ALASSO ..... ۹۱

شکل ۱۱-۲ نمودار تغییرات مقادیر پیش‌بینی شده  $pIC_{50}$  به‌وسیله مدل ALASSO-LM-ANN در شرایط بهینه در مقابل مقادیر تجربی برای داده‌های مجموعه آزمون ..... ۹۵.

شکل ۱۲-۲ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌ها بر اساس تکنیک LOO در مقابل مقادیر تجربی ..... ۹۹.

شکل ۱۳-۲ نمودار باقی‌مانده‌های حاصل از پیش‌بینی فعالیت دارویی ترکیبات با استفاده از تکنیک LOO و مقادیر تجربی برحسب مقادیر تجربی ..... ۹۹.

شکل ۱۴-۲ دامنه کاربرد مدل ALASSO-LM-ANN، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $3 \pm \sigma$  و  $h^*$  است. .... ۱۰۲.

شکل ۱۵-۲ نمودار مقادیر  $R^2$  به‌دست آمده در آزمون Y-تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای Y-تصادفی و پیش‌بینی فعالیت ترکیبات مجموعه آزمون به‌وسیله مدل ALASSO-LM-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه ..... ۱۰۴.

شکل ۱۶-۲ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های ایدز ..... ۱۱۷.

شکل ۱۷-۲ نمودار مقادیر VIF توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های ایدز ..... ۱۱۷.

شکل ۱۸-۲ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های سرطان کارسینوم کولورکتال ..... ۱۱۸.

شکل ۱۹-۲ نمودار مقادیر VIF توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های سرطان کارسینوم کولورکتال ..... ۱۱۸.

شکل ۲۰-۲ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های سرطان ریه ..... ۱۱۹.

شکل ۲۱-۲ نمودار مقادیر VIF توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های سرطان ریه ..... ۱۱۹.

شکل ۲۲-۲ نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی برای داده‌های ضد ایدز مجموعه آزمون ..... ۱۲۵.

شکل ۲۳-۲ نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی برای داده‌های ضد سرطان کارسینوم کولورکتال مجموعه آزمون ..... ۱۲۵.

شکل ۲۴-۲ نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی برای داده‌های ضد سرطان ریه مجموعه آزمون

۱۲۶.....

شکل ۲-۲۵ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌های ضد ایدز بر اساس تکنیک LOO در مقابل مقادیر تجربی

۱۳۱.....

شکل ۲-۲۶ نمودار باقی‌مانده‌های پیش‌بینی شده داده‌های ضد ایدز با استفاده از تکنیک LOO برحسب مقادیر تجربی ۱۳۱

شکل ۲-۲۷ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌های ضد سرطان کارسینوم کولورکتال بر اساس تکنیک LOO

در مقابل مقادیر تجربی ..... ۱۳۲.....

شکل ۲-۲۸ نمودار باقی‌مانده‌های پیش‌بینی شده داده‌های ضد سرطان کارسینوم کولورکتال با استفاده از تکنیک LOO

برحسب مقادیر تجربی ..... ۱۳۲.....

شکل ۲-۲۹ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌های ضد سرطان ریه بر اساس تکنیک LOO در مقابل مقادیر

تجربی ..... ۱۳۳.....

شکل ۲-۳۰ نمودار باقی‌مانده‌های پیش‌بینی شده داده‌های ضد سرطان ریه با استفاده از تکنیک LOO برحسب مقادیر تجربی

..... ۱۳۳.....

شکل ۲-۳۱ دامنه کاربرد مدل LAD-LASSO-LM-ANN برای مجموعه داده‌های ضد ایدز، خطوط نقطه چین افقی و

عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $3 \pm \sigma$  و  $h^*$  است. .... ۱۳۷.....

شکل ۲-۳۲ دامنه کاربرد مدل LAD-LASSO-LM-ANN برای مجموعه داده‌های ضد سرطان کارسینوم کولورکتال، خطوط

نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $3 \pm \sigma$  و  $h^*$  است. .... ۱۳۷.....

شکل ۲-۳۳ دامنه کاربرد مدل LAD-LASSO-LM-ANN برای مجموعه داده‌های ضد سرطان ریه، خطوط نقطه چین افقی

و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $3 \pm \sigma$  و  $h^*$  است. .... ۱۳۸.....

شکل ۲-۳۴ نمودار مقادیر  $R^2$  به دست آمده در آزمون Y-تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای Y-تصادفی و

پیش‌بینی فعالیت ترکیبات ضد ایدز مجموعه آزمون به‌وسیله مدل LAD-LASSO-LM-ANN با استفاده از پاسخ تصادفی

شده در شرایط بهینه ..... ۱۳۹.....

شکل ۲-۳۵ نمودار مقادیر  $R^2$  به دست آمده در آزمون Y-تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای Y-تصادفی و

پیش‌بینی فعالیت ترکیبات ضد سرطان کارسینوم کولورکتال مجموعه آزمون به‌وسیله مدل LAD-LASSO-LM-ANN با

استفاده از پاسخ تصادفی شده در شرایط بهینه ..... ۱۴۰.....

شکل ۲-۳۶ نمودار مقادیر  $R^2$  به دست آمده در آزمون  $Y$ -تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای  $Y$ -تصادفی و پیش‌بینی فعالیت ترکیبات ضد سرطان ریه مجموعه آزمون به وسیله مدل LAD-LASSO-LM-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه ..... ۱۴۰

شکل ۲-۳۷ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب SCAD برای مجموعه داده A. ۱۵۲

شکل ۲-۳۸ نمودار مقادیر VIF توصیف‌کننده‌های منتخب SCAD برای مجموعه داده A ..... ۱۵۲

شکل ۲-۳۹ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب SCAD برای مجموعه داده B. ۱۵۳

شکل ۲-۴۰ نمودار مقادیر VIF توصیف‌کننده‌های منتخب SCAD برای مجموعه داده B ..... ۱۵۳

شکل ۲-۴۱ نمودار تغییرات مقادیر پیش‌بینی شده RI به وسیله مدل SCAD-LM-ANN در شرایط بهینه در مقابل مقادیر تجربی برای داده‌های مجموعه آزمون مجموعه A ..... ۱۵۹

شکل ۲-۴۲ نمودار تغییرات مقادیر پیش‌بینی شده RI به وسیله مدل SCAD-BR-ANN در شرایط بهینه در مقابل مقادیر تجربی برای داده‌های مجموعه آزمون مجموعه B ..... ۱۵۹

شکل ۲-۴۳ نمودار تغییرات مقادیر پیش‌بینی شده RI همه داده‌های مجموعه A بر اساس تکنیک LOO در مقابل مقادیر تجربی ..... ۱۶۴

شکل ۲-۴۴ نمودار باقی‌مانده‌های پیش‌بینی شده RI همه داده‌های مجموعه A با استفاده از تکنیک LOO بر حسب مقادیر تجربی ..... ۱۶۴

شکل ۲-۴۵ نمودار تغییرات مقادیر پیش‌بینی شده RI همه داده‌های مجموعه B بر اساس تکنیک LOO در مقابل مقادیر تجربی ..... ۱۶۵

شکل ۲-۴۶ نمودار باقی‌مانده‌های پیش‌بینی شده RI همه داده‌های مجموعه B با استفاده از تکنیک LOO بر حسب مقادیر تجربی ..... ۱۶۵

شکل ۲-۴۷ دامنه کاربرد مدل SCAD-LM-ANN برای مجموعه A، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $3\pm\sigma$  و  $h^*$  است. ..... ۱۷۰

شکل ۲-۴۸ دامنه کاربرد مدل SCAD-BR-ANN برای مجموعه B، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $3\pm\sigma$  و  $h^*$  است. ..... ۱۷۰

شکل ۲-۴۹ نمودار مقادیر  $R^2$  به دست آمده در آزمون  $Y$ -تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای  $Y$ -تصادفی و



پیش‌بینی بیش‌بینی RI ترکیبات آزمون مجموعه A به‌وسیله مدل SCAD-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه ..... ۱۷۲

شکل ۲-۵۰ نمودار مقادیر  $R^2$  به‌دست آمده در آزمون Y-تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای Y-تصادفی و پیش‌بینی بیش‌بینی RI ترکیبات آزمون مجموعه B به‌وسیله مدل SCAD-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه ..... ۱۷۲

شکل ۳-۱ نمودار سهم مشارکت توصیف‌کننده‌ها در مدل SCAD-LM-ANN ..... ۱۷۶

شکل ۳-۲ ساختار کریستالوگرافی 3DLG [۱۹۸] (منطقه نقطه چین نشان‌دهنده لیگاند کریستالوگرافی و مابقی زنجیره‌های اسید آمینه‌ای است) ..... ۱۸۰

شکل ۳-۳ بررسی برهم‌کنش ترکیبات NC1، ۴۴، ۵۶، NC2 و NC3 با گیرنده ..... ۱۸۴

شکل ۳-۴ بررسی برهم‌کنش ترکیبات NC3، NC4، NC7 و NC8 با گیرنده ..... ۱۸۵

شکل ۳-۵ نمودار سهم مشارکت توصیف‌کننده‌ها در مدل ALASSO-LM-ANN ..... ۱۸۹

شکل ۳-۶ ساختار کریستالوگرافی 6LU7 [۱۹۸] (منطقه نقطه چین نشان‌دهنده لیگاند کریستالوگرافی و مابقی زنجیره‌های اسید آمینه‌ای است) ..... ۱۹۵

شکل ۳-۷ برهم‌کنش ترکیبات فعال (۷۴) و کم‌فعال (۷۲) موجود در مجموعه داده‌ها با اسید آمینه‌های کلیدی ..... ۱۹۸

شکل ۳-۸ برهم‌کنش ترکیبات پیشنهادی (NC1 و NC3) با اسید آمینه‌های کلیدی ..... ۱۹۹

شکل ۳-۹ برهم‌کنش ترکیبات پیشنهادی (NC4 و NC7) با اسید آمینه‌های کلیدی ..... ۲۰۰

شکل ۳-۱۰ برهم‌کنش ترکیبات پیشنهادی (NC8 و NC26) با اسید آمینه‌های کلیدی ..... ۲۰۱

شکل ۳-۱۱ برهم‌کنش ترکیبات پیشنهادی (NC28 و NC30) با اسید آمینه‌های کلیدی ..... ۲۰۲

شکل ۳-۱۲ نمودار سهم مشارکت توصیف‌کننده‌ها در مدل LAD-LASSO-LM-ANN برای بازدارنده‌های ضد ایدز .. ۲۱۰

شکل ۳-۱۳ نمودار سهم مشارکت توصیف‌کننده‌ها در مدل LAD-LASSO-LM-ANN برای بازدارنده‌های سرطان کارسینوم کولورکتال ..... ۲۱۱

شکل ۳-۱۴ نمودار سهم مشارکت توصیف‌کننده‌ها در مدل LAD-LASSO-LM-ANN برای بازدارنده‌های سرطان ریه ..... ۲۱۱

شکل ۳-۱۵ ساختار کریستالوگرافی 3M8Q [۲۱۵] (منطقه نقطه چین نشان‌دهنده لیگاند کریستالوگرافی و مابقی زنجیره‌های

- اسید آمینه‌ای است)..... ۲۱۶.....
- شکل ۱۶-۳ برهم‌کنش ترکیب فعال (۲۰) موجود در مجموعه داده‌های ضد ایدز با اسید آمینه‌های کلیدی ..... ۲۱۸.....
- شکل ۱۷-۳ برهم‌کنش ترکیب کم فعال (۴۰) موجود در مجموعه داده‌های ضد ایدز با اسید آمینه‌های کلیدی ..... ۲۱۸.....
- شکل ۱۸-۳ برهم‌کنش ترکیب کم فعال (۳۶) موجود در مجموعه داده‌های ضد ایدز با اسید آمینه‌های کلیدی ..... ۲۱۹.....
- شکل ۱۹-۳ برهم‌کنش ترکیب پیشنهادی NC1 با اسید آمینه‌های کلیدی ..... ۲۱۹.....
- شکل ۲۰-۳ برهم‌کنش ترکیب پیشنهادی NC2 با اسید آمینه‌های کلیدی ..... ۲۲۰.....
- شکل ۲۱-۳ برهم‌کنش ترکیب پیشنهادی NC3 با اسید آمینه‌های کلیدی ..... ۲۲۰.....
- شکل ۲۲-۳ برهم‌کنش ترکیب پیشنهادی NC4 با اسید آمینه‌های کلیدی ..... ۲۲۱.....
- شکل ۲۳-۳ برهم‌کنش ترکیب پیشنهادی NC5 با اسید آمینه‌های کلیدی ..... ۲۲۱.....
- شکل ۲۴-۳ برهم‌کنش ترکیب پیشنهادی NC6 با اسید آمینه‌های کلیدی ..... ۲۲۲.....
- شکل ۲۵-۳ برهم‌کنش ترکیب پیشنهادی NC7 با اسید آمینه‌های کلیدی ..... ۲۲۲.....
- شکل ۲۶-۳ ساختار کریستالوگرافی 3HHM [۲۱۸] (منطقه نقطه چین نشان‌دهنده لیگاند کریستالوگرافی و مابقی زنجیره‌های اسید آمینه‌ای است)..... ۲۳۱.....
- شکل ۲۷-۳ برهم‌کنش ترکیب نسبتاً فعال (۶۸) موجود در مجموعه داده‌های ضد سرطان کاسینوم کولورکتال با اسید آمینه‌های کلیدی ..... ۲۳۲.....
- شکل ۲۸-۳ برهم‌کنش ترکیب کم فعال (۲۷) موجود در مجموعه داده‌های ضد سرطان کاسینوم کولورکتال با اسید آمینه‌های کلیدی ..... ۲۳۳.....
- شکل ۲۹-۳ برهم‌کنش ترکیب پیشنهادی NC1 با اسید آمینه‌های کلیدی ..... ۲۳۳.....
- شکل ۳۰-۳ برهم‌کنش ترکیب پیشنهادی NC2 با اسید آمینه‌های کلیدی ..... ۲۳۴.....
- شکل ۳۱-۳ برهم‌کنش ترکیب پیشنهادی NC3 با اسید آمینه‌های کلیدی ..... ۲۳۴.....
- شکل ۳۲-۳ برهم‌کنش ترکیب پیشنهادی NC4 با اسید آمینه‌های کلیدی ..... ۲۳۵.....
- شکل ۳۳-۳ نمودار سهم مشارکت توصیف‌کننده‌ها در مدل SCAD-ANN برای مجموعه A ..... ۲۳۸.....
- شکل ۳۴-۳ نمودار سهم مشارکت توصیف‌کننده‌ها در مدل SCAD-ANN برای مجموعه B ..... ۲۳۹.....

## فهرست جدول‌ها:

جدول ۱-۱ پارامترهای آماری .....	۲۹
جدول ۱-۲ ساختار ترکیبات شیمیایی به همراه مقادیر واقعی و پیش‌بینی شده $pEC_{50}$ .....	۶۳
جدول ۲-۲ توصیف‌کننده‌های منتخب SCAD .....	۶۶
جدول ۳-۲ ساختارهای شبکه‌های توسعه یافته با توصیف‌کننده‌های منتخب SCAD با کمترین MSE مجموعه ارزیابی .	۷۰
جدول ۴-۲ نتایج حاصل از ارزیابی مدل SCAD-ANN با استفاده از مجموعه آزمون .....	۷۲
جدول ۵-۲ نتایج حاصل از ارزیابی مدل SCAD-LM-ANN به روش رد مرحله‌ای تک تک برای کل داده‌ها .....	۷۵
جدول ۶-۲ پارامترهای آماری محاسبه شده برای مجموعه آزمون و داده‌های پیش‌بینی شده با تکنیک LOO برای مدل SCAD-LM-ANN .....	۷۸
جدول ۷-۲ مجموعه داده‌ها به همراه مقادیر واقعی و پیش‌بینی شده $pIC_{50}$ .....	۸۶
جدول ۸-۲ توصیف‌کننده‌های منتخب ALASSO .....	۹۰
جدول ۹-۲ ساختارهای شبکه‌های توسعه یافته با توصیف‌کننده‌های منتخب ALASSO با کمترین MSE مجموعه ارزیابی .....	۹۳
جدول ۱۰-۲ نتایج حاصل از ارزیابی مدل ALASSO-LM-ANN با استفاده از مجموعه آزمون .....	۹۵
جدول ۱۱-۲ نتایج حاصل از ارزیابی مدل ALASSO-LM-ANN به روش رد مرحله‌ای تک تک برای کل داده‌ها .....	۹۷
جدول ۱۲-۲ پارامترهای آماری محاسبه شده برای مجموعه آزمون و داده‌های پیش‌بینی شده با تکنیک LOO برای مدل ALASSO-LM-ANN .....	۱۰۱
جدول ۱۳-۲ مجموعه داده‌های ضد ایدز به همراه مقادیر واقعی و پیش‌بینی شده $pEC_{50}$ .....	۱۰۸
جدول ۱۴-۲ مجموعه داده‌های ضد سرطان کارسینوم کولورکتال به همراه مقادیر واقعی و پیش‌بینی شده $pIC_{50}$ .....	۱۱۰
جدول ۱۵-۲ مجموعه داده‌های ضد سرطان ریه به همراه مقادیر واقعی و پیش‌بینی شده $pIC_{50}$ .....	۱۱۲
جدول ۱۶-۲ توصیف‌کننده‌های منتخب LAD-LASSO برای هر سه مجموعه داده‌ها .....	۱۱۶
جدول ۱۷-۲ ساختارهای شبکه‌های توسعه یافته با توصیف‌کننده‌های منتخب LAD_LASSO با کمترین MSE مجموعه	

ارزیابی.....	۱۲۱
جدول ۱۸-۲ نتایج حاصل از ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از مجموعه آزمون.....	۱۲۴
جدول ۱۹-۲ نتایج حاصل از ارزیابی مدل LAD-LASSO-LM-ANN با تکنیک LOO برای کل مجموعه داده‌های ضد ایدز.....	۱۲۸
جدول ۲۰-۲ نتایج حاصل از ارزیابی مدل LAD-LASSO-LM-ANN به روش رد مرحله‌ای تک تک برای کل مجموعه داده‌های ضد سرطان کارسینوم کولورکتال.....	۱۲۹
جدول ۲۱-۲ نتایج حاصل از ارزیابی مدل LAD-LASSO-LM-ANN به روش رد مرحله‌ای تک تک برای کل مجموعه داده‌های ضد سرطان ریه.....	۱۳۰
جدول ۲۲-۲ پارامترهای آماری محاسبه شده برای مجموعه آزمون و داده‌های پیش‌بینی شده با تکنیک LOO برای مدل LAD-LASSO-LM-ANN هر سه مجموعه از داده‌ها.....	۱۳۵
جدول ۲۳-۲ ساختار شیمیایی و SMILES مربوط به ترکیبات مورد مطالعه مجموعه داده‌های A به همراه شاخص بازداری.....	۱۴۵
جدول ۲۴-۲ ساختار شیمیایی و SMILES مربوط به ترکیبات مورد مطالعه مجموعه داده‌های B به همراه شاخص بازداری.....	۱۴۸
جدول ۲۵-۲ توصیف‌کننده‌های منتخب SCAD برای مجموعه داده‌های A و B.....	۱۵۱
جدول ۲۶-۲ ساختارهای شبکه‌های توسعه یافته با توصیف‌کننده‌های منتخب SCAD با کمترین MSE مجموعه ارزیابی هر دو مجموعه داده A و B.....	۱۵۶
جدول ۲۷-۲ نتایج حاصل از ارزیابی مدل SCAD-ANN با استفاده از مجموعه آزمون.....	۱۵۸
جدول ۲۸-۲ نتایج حاصل از ارزیابی مدل SCAD-ANN با تکنیک LOO برای کل داده‌های مجموعه A.....	۱۶۱
جدول ۲۹-۲ نتایج حاصل از ارزیابی مدل SCAD-ANN با تکنیک LOO برای کل داده‌های مجموعه B.....	۱۶۳
جدول ۳۰-۲ پارامترهای آماری محاسبه شده برای مجموعه آزمون و داده‌های پیش‌بینی شده با تکنیک LOO برای مدل SCAD-ANN هر دو مجموعه از داده‌ها.....	۱۶۸
جدول ۱-۳ پارامترهای PK محاسبه شده برای ترکیبات مورد مطالعه و ترکیبات پیشنهادی.....	۱۸۶
جدول ۲-۳ پارامترهای PK محاسبه شده برای ترکیبات مورد مطالعه و ترکیبات پیشنهادی.....	۲۰۳

- جدول ۳-۳ پارامترهای PK محاسبه شده برای ترکیبات مورد مطالعه ضد ایدز و ترکیبات پیشنهادی ..... ۲۲۳
- جدول ۴-۳ پارامترهای PK محاسبه شده برای ترکیبات مورد مطالعه ضد سرطان و ترکیبات پیشنهادی ..... ۲۳۶
- جدول ۵-۳ مقایسه پارامترهای آماری محاسبه شده برای مجموعه آزمون مدل برتر SCAD-LM-ANN با مدل SCAD برای مشتقات استانیلید/ استامید بهعنوان بازدارنده‌های ایدز ..... ۲۴۷
- جدول ۶-۳ مقایسه پارامترهای آماری محاسبه شده برای مجموعه آزمون مدل برتر ALASSO-LM-ANN با مدل ALASSO برای مشتقات (3-chymotrypsin like protease (3CLPro) بهعنوان بازدارنده‌های SARS-COV-2 ..... ۲۴۸
- جدول ۷-۳ مقایسه پارامترهای آماری محاسبه شده برای مجموعه آزمون مدل برتر LAD-LASSO-ANN با مدل LAD- LASSO برای هر سه مجموعه داده‌های متفاوت ..... ۲۴۹
- جدول ۸-۳ مقایسه پارامترهای آماری محاسبه شده برای مجموعه آزمون مدل برتر SCAD -ANN با مدل SCAD برای مجموعه داده A و B ..... ۲۵۰

## فهرست رابطه‌ها:

۱۴.....	رابطه ۱-۱.....
۱۴.....	رابطه ۲-۱.....
۱۷.....	رابطه ۳-۱.....
۱۷.....	رابطه ۴-۱.....
۱۷.....	رابطه ۵-۱.....
۱۸.....	رابطه ۶-۱.....
۱۸.....	رابطه ۷-۱.....
۱۸.....	رابطه ۸-۱.....
۱۹.....	رابطه ۹-۱.....
۲۰.....	رابطه ۱۰-۱.....
۲۳.....	رابطه ۱۱-۱.....
۲۶.....	رابطه ۱۲-۱.....
۳۰.....	رابطه ۱۳-۱.....
۳۰.....	رابطه ۱۴-۱.....
۳۲.....	رابطه ۱۵-۱.....
۱۷۶.....	رابطه ۱-۳.....
۱۸۹.....	رابطه ۲-۳.....
۲۳۹.....	رابطه ۳-۳.....
۲۳۹.....	رابطه ۴-۳.....
۲۴۰.....	رابطه ۵-۳.....

## فهرست کلمات اختصاری

QSAR (Quantitative structure- activity relationships)  
QSPR (Quantitative structure- properties relationships)  
CADD (Computer Aided Drug Design)  
RI (Retention index)  
OECD (The Organization for Economic Co-operation and Development)  
KS (Kennard-Stone)  
LASSO (Least absolute deviation- Least absolute shrinkage and selection operator)  
LAD (Least Absolute Deviations)  
SCAD( Smoothly Clipped Absolute Deviation)  
VIF (Variance inflation factor)  
PCR (Principal component regression)  
PLS (Partial least square)  
ANN (Artificial neural network)  
LM (Levenberg –Marquardt)  
BR (Bayesian regularization)  
LOO (Leave-one-out)  
AD (Applicability Domain)  
CV (Cross validation)  
RMSD (Root-mean-square deviation)  
SARS (Severe acute respiratory syndrome coronavirus)  
MERS (Middle East respiratory syndrome coronavirus)  
LAD- LASSO (Least absolute deviation (LAD)- least absolute shrinkage and selection operator (LASSO))  
VOCs (Volatile organic compounds)  
OLS (Ordinary least squares)  
RMSE (Root mean square error)  
LASSO (Least Absolute Shrinkage and Selection Operator)  
ALASSO (Adaptive least Absolute Shrinkage and Selection Operator)  
MCP (Minimax concave penalty)  
SPSS (Statistical Package for Social Science)  
SMILES (Simplified molecular-input line-entry system)  
SIS (Sure Independence Screening)  
ISIS (Iterative Sure Independence Screening)





مقدمه



## ۱- مقدمه

یکی از مشکلاتی که انسان همیشه با آن مواجه است، مقابله و درمان انواع بیماری‌هایی است که زندگی بشریت را با خطر روبرو می‌کند و همواره یکی از مهم‌ترین مسائل چالش‌برانگیز دانشمندان، توسعه داروهای مؤثر برای رفع و یا کاهش اثرات بیماری‌ها می‌باشد. بیماری‌هایی مانند ایدز، انواع سرطان و مقاومت ویروس‌ها در برابر انواع داروها و جهش‌های متنوعی که ایجاد می‌شود، منجر به تلاش مستمر محققین برای یافتن داروهای مؤثر و کارآمد در مواجهه با بیماری‌ها می‌شود.

توسعه دارو در گذشته به روش‌های آزمون و خطا انجام می‌شده است که به دلایل صرف زمان و هزینه مضاعف همواره دردسر ساز بوده است. علاوه بر این دانشمندان با توجه به مسائلی همچون عدم آگاهی اولیه از امکان سنتز و نتایج آزمایشگاهی دال بر یک سنتز موفق، عدم اطلاع از میزان فعالیت دارویی ترکیبات بالقوه و در دسترس نبودن امکانات انجام آنالیزهای متفاوت، همیشه در مسیر توسعه داروها با مشکلات عدیده‌ای روبرو می‌شوند. از جمله مهم‌ترین اهداف محققین و شیمیدان‌ها می‌توان به پیش‌بینی فعالیت دارویی ترکیبات قبل از سنتز اشاره کرد. زیرا با تحقق این هدف، در بخش عمده‌ای از هزینه‌ها، زمان و به‌کارگیری نیروی انسانی صرفه‌جویی خواهد شد.

علاوه بر مشکلات اشاره شده در توسعه داروها، پیش‌بینی خواص مربوط به ترکیبات شیمیایی نیز همواره مورد توجه محققین و شیمیدان‌ها بوده است. به‌طوری‌که با توجه به نبود امکانات کامل و جامع آزمایشگاهی و به‌مخاطره نیافتن طبیعت، تعیین مقدار برخی از خواص شیمیایی و فیزیکی ترکیبات نیز دارای اهمیت منحصر به فردی می‌باشد. بنابراین استفاده از روش‌های تئوری و شیمی محاسباتی توانسته است دانشمندان را در تحقق این هدف همراهی کند، تا بدون انجام آزمایش‌های متعدد فعالیت دارویی ترکیبات قبل از سنتز پیش‌بینی شود. علاوه بر این با به‌کارگیری شیمی محاسباتی، پیش‌بینی خواص فیزیکوشیمیایی ترکیبات نیز امکان‌پذیر می‌باشد.

شیمی نظری محاسباتی اساساً به محاسبه عددی ساختارهای مولکولی و برهمکنش‌های مولکولی می‌پردازد. اصطلاح شیمی محاسباتی معمولاً زمانی استفاده می‌شود که یک روش ریاضی به اندازه کافی توسعه یافته باشد که بتوان آن را برای پیاده‌سازی بر روی رایانه، خودکار کرد. شیمی محاسباتی استفاده از مهارت‌های شیمی، ریاضی و رایانه برای حل مسائل جالب شیمیایی است. از رایانه‌ها برای تولید اطلاعاتی مانند ویژگی‌های مولکول‌ها یا شبیه‌سازی برهم‌کنش ترکیبات استفاده می‌شود. شیمی محاسباتی، روشی مفید برای بررسی مواد شیمیایی است که یافتن آن‌ها بسیار دشوار یا خرید آن‌ها بسیار گران است. همچنین به شیمیدانان کمک می‌کند قبل از اجرای آزمایش‌های واقعی، برخی از ویژگی‌های ترکیبات مورد مطالعه را پیش‌بینی کنند تا بتوانند به‌طور کارآمدی برای انجام آزمایش‌ها آماده شوند. از طرفی هزینه‌های آزمایشگاهی و به‌کارگیری نیروی انسانی و به‌پیروی از آن خطرات قرارگیری در معرض مواد شیمیایی نیز کاهش می‌یابد. بدون شک، امروزه بدون شبیه‌سازی مولکولی و بررسی‌های تئوری اولیه، هیچ سنتز آزمایشگاهی و آنالیز دستگاهی هدفمند و به‌صرفه نخواهد بود [۱].

## ۱-۲ کمومتریکس

کمومتریکس<sup>۲</sup> با به‌کارگیری روش‌های داده‌محور به استخراج داده‌ها از سیستم‌های شیمیایی می‌پردازد. در سال ۱۹۷۰ با توجه به پیشرفت الکترونیک و به‌کارگیری فزاینده کامپیوترها در علوم مختلف، این علم شناخته شد. اصطلاح کمومتریکس برای اولین بار در سال ۱۹۷۱، توسط سوانت ولد<sup>۳</sup> بیان شد و انجمن بین‌المللی کمومتریکس توسط بروس کوالسکی<sup>۴</sup> استاد شیمی تجزیه دانشگاه واشنگتن و ولد، تشکیل داده شد [۲]. علم کمومتریکس یک دانش بین‌رشته‌ای است، به‌طوری‌که با استفاده از علم شیمی محاسباتی، آمار کاربردی و کامپیوتر به تحلیل سیستم‌های متفاوت در شیمی، بیوفیزیک، بیوشیمی، زیست‌شناسی و

---

Computational theoretical chemistry

<sup>۲</sup>Chemometrics

<sup>۳</sup>Svante Wold

<sup>۴</sup>Bruce Kowalski

مهندسی شیمی می‌پردازد. به عبارت دیگر، کمومتریکس به معنای استفاده از روش‌های محاسباتی در اندازه‌گیری داده‌های شیمیایی می‌باشد. شیمیدان‌های تجزیه با به‌کارگیری تکنیک‌های متفاوت کمومتریکس، به پیشرفت‌های متنوعی در علوم دستگاهی، طراحی و انتخاب فرایندها و بهبود در طراحی و توسعه داروها به کمک روش‌های مدل‌سازی قدرتمند دست یافته‌اند.

بنابراین با توجه به کاربردهای ویژه تکنیک‌های علم کمومتریکس، می‌توان از آن برای توسعه ارتباط کمی ساختار- فعالیت<sup>۱</sup> (QSAR) و یا ارتباط کمی ساختار- خاصیت<sup>۲</sup> (QSPR) استفاده کرد.

### ۱-۳ رابطه کمی ساختار- فعالیت (QSAR)

رابطه کمی ساختار- فعالیت (QSAR) یک مدل ریاضی و محاسباتی برای آشکار کردن روابط بین فعالیت‌های دارویی و ویژگی‌های ساختاری ترکیبات شیمیایی است. QSAR به‌عنوان یک روش طراحی دارو بیش از ۵۰ سال پیش توسط هانش و فوجیتا<sup>۳</sup> توسعه یافت [۳]. از آن زمان تاکنون، QSAR یک روش کارآمد برای ساخت مدل‌های ریاضی است که تلاش می‌کند تا با استفاده از تکنیک‌های رگرسیونی، یک همبستگی آماری معنادار بین ساختار شیمیایی ترکیبات و فعالیت دارویی ( $K_i$ ,  $pEC_{50}$ ,  $pIC_{50}$ ، و غیره) پیدا کند [۴]. اصل اساسی این روش به‌صورتی است که تغییرات در ویژگی‌های ساختاری (خواص فیزیکی و شیمیایی) باعث تغییرات متفاوتی در فعالیت‌های دارویی می‌شود [۵]. مدل‌سازی QSAR تعداد زیادی از مواد شیمیایی را از نظر فعالیت‌های دارویی مورد نظر اولویت‌دهی می‌کند. این روش یک فرایند شبیه‌سازی رایانه‌ای<sup>۴</sup> است که با پیش‌بینی فعالیت دارویی ترکیبات، تعداد مواد شیمیایی کاندید برای انجام آزمایش را به میزان قابل توجهی کاهش می‌دهد [۶, ۷].

<sup>۱</sup>Quantitative structure- activity relationships

<sup>۲</sup>Quantitative structure- properties relationships

<sup>۳</sup>Hansch and Fujita

<sup>۴</sup>In silico

در توسعه مدل‌های QSAR، باید از مجموعه‌ای از داده‌های همگن استفاده شود، به این معنی که فعالیت دارویی آن‌ها، در شرایط آزمایشگاهی مشابه و با روش اندازه‌گیری یکسان اندازه‌گیری شده باشد، تا مدل نهایی از اعتبار قابل قبولی برخوردار باشد [۸]. استفاده از مجموعه داده‌های همگن منجر به انتخاب توصیف‌کننده‌های کم‌تر می‌شود و از این‌رو مدل‌های QSAR با تفسیرپذیری بالاتری ایجاد می‌شود. توصیف‌کننده‌های مولکولی مقادیر عددی و بیانگر ویژگی‌های فیزیکوشیمیایی ترکیبات مورد مطالعه هستند. توصیف‌کننده‌های مولکولی به دو صورت محاسباتی و یا تجربی استخراج می‌شوند و ساختار مولکولی را کمی می‌کنند و نشان‌دهنده اطلاعات مربوط به ساختار شیمیایی هستند. از ساده‌ترین توصیف‌کننده‌ها، می‌توان به تعداد و انواع اتم‌ها یا نوع پیوندهای شیمیایی موجود در ساختار مولکول اشاره کرد. توصیف‌کننده‌های مولکولی دقیق‌تر را نیز می‌توان از طریق نظریه‌های مختلف، مانند مکانیک کوانتومی و معادلات پیچیده‌تر ریاضی استخراج کرد [۹].

توسعه مدل‌های QSAR برای ترکیبات شیمیایی سنتز شده این امکان را به شیمیدان‌ها می‌دهد تا ترکیباتی با فعالیت دارویی بهتر پیشنهاد دهند. امروزه شیوع بیماری‌های همه‌گیر، نبود درمان و داروهای مؤثر برای برخی از بیماری‌ها، زندگی بشریت را در مواجهه با این بیماری‌ها با مخاطره همراه نموده است. بنابراین دانشمندان در تلاش هستند تا در راستای مهم‌ترین دغدغه بشریت که حفظ و ارتقا سلامتی و درمان بیماری‌های انسانی است، گام رو به رشدی بردارند. از جمله بیماری‌های چالش برانگیز می‌توان به ایدز، سرطان، آلزایمر، بیماری‌های همه‌گیر جدید همچون کووید-۱۹ اشاره کرد. البته مقاومت ویروس‌ها در برابر داروها و جهش‌های ویروسی همگی منجر به تلاش دانشمندان در پیشبرد تحقیقات بهبود طراحی و توسعه دارو شده است. امروزه با توجه به پیشرفت علوم الکترونیک و وجود کامپیوترها و نرم‌افزارهای محاسباتی قدرتمند، طراحی دارو به روش منطقی و مبتنی بر اصول محاسباتی به‌خوبی و با موفقیت جایگزین مناسبی برای طراحی داروی سنتی می‌باشد. بنابراین با توجه به مزایای یادشده، طراحی دارو به کمک

کامپیوتر<sup>۱</sup> به دلیل صرفه‌جویی در زمان و هزینه‌های پژوهشی و کاهش نیروی انسانی لازم، ابزار مناسبی در پیشبرد فرایند توسعه دارو بوده است و ذهن دانشمندان طراحی دارو را به خود معطوف نموده است.

## ۱-۴ رابطه کمی ساختار - خاصیت (QSPR)

رابطه کمی ساختار - خاصیت (QSPR) یک رابطه ریاضی بین ویژگی‌های ساختاری و خاصیت مربوط به مجموعه‌ای از مواد شیمیایی را توصیف می‌کند. طبیعتاً استفاده از چنین روابط ریاضی برای پیش‌بینی ویژگی هدف انواع مواد شیمیایی قبل از اندازه‌گیری‌های تجربی پرهزینه، فشرده و زمان‌بر مورد توجه محققین شیمیدان بوده است. استفاده از مدل‌های QSPR برای غربالگری پایگاه‌های اطلاعاتی شیمیایی قبل از سنتز و اندازه‌گیری تجربی آن‌ها، برای تولیدکنندگان مواد شیمیایی، شرکت‌های دارویی و سازمان‌های دولتی، به‌ویژه در زمان‌هایی که با کمبود منابع انسانی، بودجه‌ای و یا امکانات آزمایشگاهی مواجه بوده‌اند، بسیار حائز اهمیت است. با توجه به روند رو به رشد پایگاه‌های ترکیبات شیمیایی سنتز شده از یک‌طرف و فشارهای قانونی و اجتماعی برای ارزیابی به‌موقع خطرات بهداشتی و زیست محیطی مواد شیمیایی از سوی دیگر، تحریم‌های سیاسی موجود و ممنوعیت خرید مواد شیمیایی اولیه و استانداردهای شیمیایی از بازار جهانی و تورم اقتصادی موجود در جوامع جهان سوم، نیاز به ارائه مدل‌های QSPR قابل اعتماد همواره ضروری است [۱۰]. QSPR با استفاده از روش‌های کمومتریکس به بررسی چگونگی تأثیر تغییرات ساختار شیمیایی مولکول بر یک خاصیت شیمیایی معین می‌پردازد. بنابراین با استفاده از شیمی محاسباتی می‌توان بینش جدیدی از نظریه‌های مبتنی بر ارائه مدل‌های QSPR را درک نمود [۱۱، ۱۲]. از مطالعات QSPR برای تخمین خواصی مانند چگالی، نقطه جوش، حلالیت، شاخص بازدارندگی کروماتوگرافی

---

<sup>۱</sup>Computer Aided Drug Design

(RI) و فشار بخار مواد شیمیایی استفاده می‌شود. استراتژی اصلی QSPR یافتن یک رابطه کمی بهینه است که با استفاده از آن بتوان خواص ترکیبات شیمیایی را پیش‌بینی کرد [۱۳].

## ۱-۵ اهمیت ساخت مدل‌های ارتباط کمی ساختار-فعالیت / ویژگی

### (QSAR/QSPR)

کیفیت، ایمنی و کارایی یک داروی راه‌یافته به بازار، با انجام آزمایش‌های مختلف مورد ارزیابی قرار می‌گیرد [۱۴]. تولید دارو فرآیندی زمان‌بر و پرهزینه است و از مرحله اولیه کشف ترکیب بالقوه دارویی تا توسعه و عرضه آن، اغلب به‌طور متوسط ۱۵ سال طول می‌کشد [۱۵] و هزینه عرضه هر دارو به بازار حدود ۹۰۰ میلیون دلار برآورد شده است [۱۶]. هزینه بالا و فرآیند طولانی به دلیل ریسک بالای شکست توسعه دارو است. تخمین زده شده است که تنها ۱۱ درصد از داروهایی که مرحله توسعه را به پایان رساندند، توسط سازمان جهانی غذا و دارو ایالت متحده تأیید شدند [۱۷]. طبق تحقیقات انجام شده، مشخص شد که ۱۰ درصد از شکست توسعه دارو به دلیل ویژگی‌های فارماکوکینتیک ضعیف بوده است، در حالی که در مرحله بالینی، ۳۰ درصد حذف ترکیب رهبر، ناشی از عدم کارایی و ۳۰ درصد دیگر ناشی از سمیت یا ایمنی بالینی است [۱۷، ۱۸]. بنابراین، پیش‌بینی این شکست‌ها قبل از مرحله بالینی به‌منظور کاهش هزینه‌های توسعه دارو مفید خواهد بود. ادعا می‌شود که با ۱۰ درصد بهبود پیش‌بینی می‌توان به‌صرفه جویی ۱۰۰ میلیون دلاری در هزینه‌های توسعه هر دارو دست یافت [۱۸]. بنابراین، روش‌های مختلفی مانند روش‌های *in vitro* یا *in vivo* یا *in silico* در اوایل مرحله توسعه دارو برای فیلتر کردن شکست‌های احتمالی استفاده می‌شوند. نمونه‌ای از روش *in silico* مدل‌های رابطه کمی ساختار-فعالیت (QSAR) است که می‌تواند برای درک عملکرد دارو، طراحی ترکیبات جدید و نمایش اطلاعات شیمیایی مناسب، مورد استفاده قرار گیرد [۱۹]-

---

<sup>1</sup>Retention index

۲۲]. بنابراین، می‌توان آزمایش‌های پرهزینه یا آنالیز خاصیت فیزیکوشیمیایی مواد خطرناک و سمی یا ترکیبات ناپایدار را با بررسی ارتباط کمی توصیف‌کننده‌های مولکولی - خاصیت شیمیایی جایگزین کرد، که این امر به‌نوبه خود می‌تواند برای پیش‌بینی پاسخ‌های هدف ترکیبات جدید استفاده شود.

### ۱-۵-۱ مراحل ساخت مدل‌های QSAR/QSPR

همان‌طور که گفته شد هدف از انجام مطالعات QSAR/QSPR، یافتن ارتباط بین توصیف‌کننده‌های ساختاری و ویژگی هدف (فعالیت دارویی و یا خاصیت فیزیکوشیمیایی) است. اما پیدا کردن این ارتباط در طی یک روش مستقیم حاصل نمی‌شود. بلکه مراحل متفاوتی برای ساخت مدل‌های QSAR/QSPR انجام می‌شود که شکل ۱-۱ نمایانگر تمامی مراحل مورد نیاز می‌باشد. همان‌طور که شکل ۱-۱ نشان می‌دهد مراحل اصلی ساخت مدل QSAR/QSPR با شماره‌های ۱ تا ۶ مشخص شده است. یک مدل QSAR/QSPR معتبر باید از ۵ اصل اساسی منتشر شده توسط سازمان همکاری اقتصادی و توسعه (OECD)، پیروی کند که به ترتیب زیر است [۲۳]:

(۱) وجود یک پاسخ هدف تعریف شده

(۲) استفاده از یک الگوریتم بدون ابهام

(۳) وجود دامنه کاربرد تعریف شده معتبر

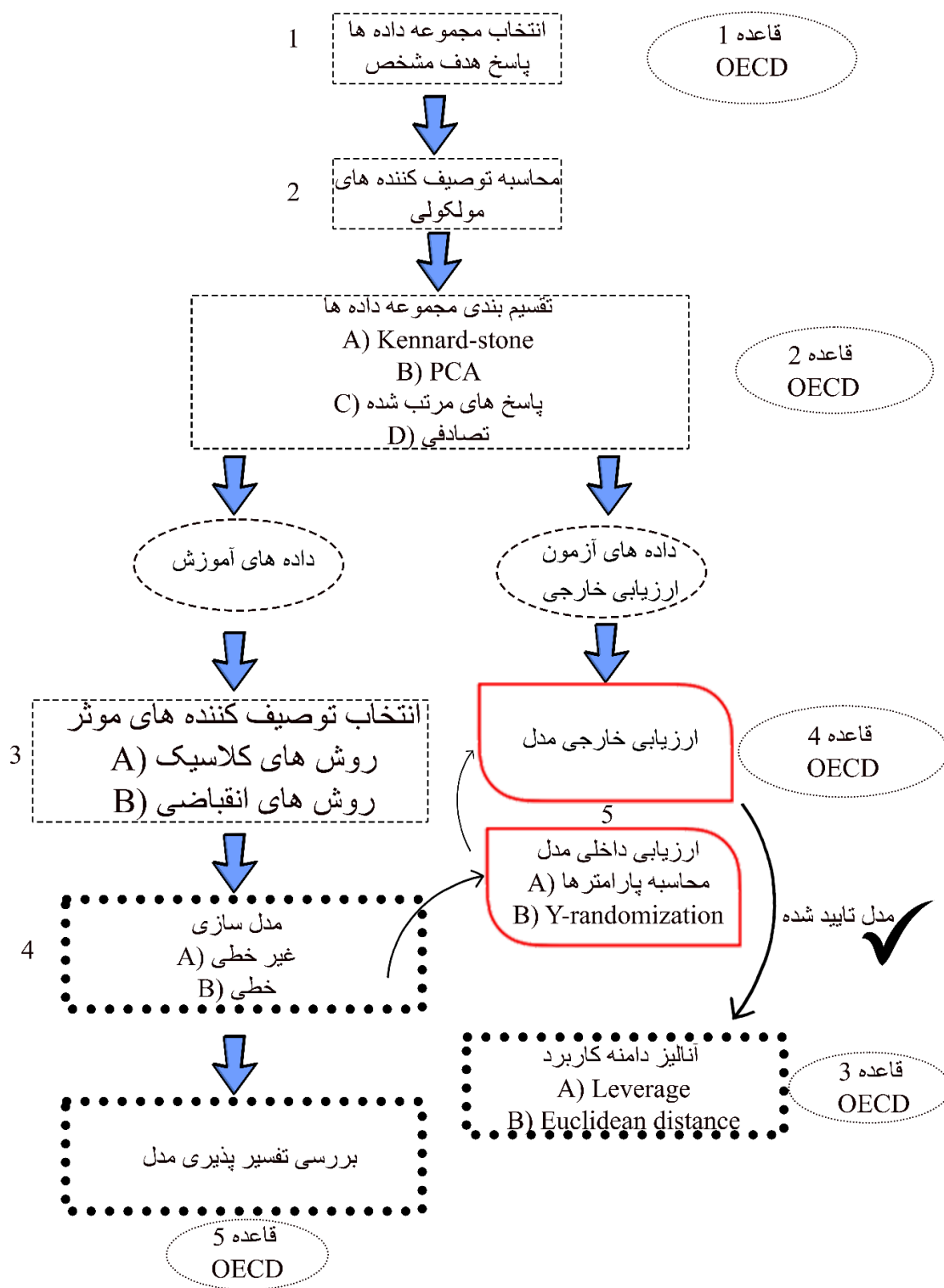
(۴) تأیید پارامترهای آماری مدل و قدرت پیش‌بینی مناسب مدل

(۵) بررسی تفسیرپذیری مدل در صورت امکان

شکل ۱-۳ جزئیات هر یک از مراحل ساخت مدل QSAR/QSPR را نشان می‌دهد. ارتباط هر کدام از بندهای قاعده OECD به مراحل ساخت مدل QSAR در شکل ۱-۳ آورده شده است.

<sup>1</sup>The Organisation for Economic Co-operation and Development





شکل ۱-۱ مراحل ساخت، توسعه و ارزیابی مدل های QSAR/QSPR

## ۱-۵-۲ جمع‌آوری و انتخاب مجموعه داده‌ها

اولین مرحله مدل‌سازی QSAR/QSPR، جمع‌آوری و انتخاب ترکیبات شیمیایی مورد نظر از منابع قابل اعتماد و در دسترس است. پاسخ هدف مورد نظر برای مدل‌سازی QSAR/QSPR، بایستی در شرایط عملی یکسان اندازه‌گیری شده باشد تا نتیجه قابل قبول‌تر و مناسب‌تری به دست آید. در مدل‌سازی QSAR/QSPR مجموعه داده‌ها با استفاده از آنالیزهای متفاوت تقسیم‌بندی از جمله کنارد-استون<sup>۱</sup> (KS) انجام می‌شود [۲۴، ۲۵]. تقسیم‌بندی KS بر اساس محاسبه فاصله اقلیدسی (Euclidean distance) است. بنابراین مجموعه داده‌ها به سه دسته آموزشی<sup>۲</sup>، ارزیابی<sup>۳</sup> و آزمون<sup>۴</sup> تقسیم می‌شود. الگوریتم KS داده‌ها را به طوری تقسیم می‌کند که مجموعه‌های تقسیم‌بندی شده مانند مجموعه ارزیابی و مجموعه آزمون نماینده مناسبی از مجموعه آموزش باشد. از این‌رو داده‌های مجموعه آموزش، ارزیابی و آزمون به ترتیب برای ساخت مدل، ارزیابی داخلی مدل و ارزیابی خارجی مدل به کار گرفته می‌شوند. لازم به ذکر است که مجموعه آزمون در هیچ یک از مراحل انتخاب متغیر و ساخت مدل‌های QSAR/QSPR شرکت نخواهد داشت.

## ۱-۵-۳ رسم و بهینه‌سازی ساختارهای شیمیایی

به منظور محاسبه صحیح توصیف‌کننده‌ها، به ساختارهای شیمیایی با حالت پایدار مولکول نیاز می‌باشد. ایجاد پایدارترین حالت ساختار شیمیایی ترکیبات با حداقل انرژی با استفاده از نرم‌افزارهای متفاوتی از جمله هایپرکم<sup>۵</sup> امکان‌پذیر است [۲۶]. بنابراین پس از رسم ساختار شیمیایی ترکیبات مورد نظر، بهینه‌سازی آن‌ها با استفاده از روش‌های متفاوتی همچون روش بهینه‌سازی نیمه تجربی (AM1) انجام

---

<sup>۱</sup>Kennard-Stone

<sup>۲</sup>Training

<sup>۳</sup>Validation

<sup>۴</sup>Test

<sup>۵</sup>Hyperchem

می‌شود. بنابراین پایدارترین حالت با کم‌ترین انرژی برای هر ساختار به دست می‌آید. ساختارهای بهینه شده به‌عنوان ورودی نرم‌افزارهای محاسبه‌ای توصیف‌کننده‌های مولکولی به کار گرفته می‌شوند.

## ۱-۵-۴ استخراج توصیف‌کننده‌ها

توصیف‌کننده‌های مولکولی، مقادیر عددی و بیانگر ویژگی‌های مربوط به ترکیبات شیمیایی هستند. توصیف‌کننده‌های مولکولی به دو دسته اصلی توصیف‌کننده‌های تجربی و محاسباتی تقسیم می‌شوند. توصیف‌کننده‌های تجربی با استفاده از پارامترهای آزمایشگاهی اندازه‌گیری می‌شوند و تنها برای ترکیبات سنتز شده و در دسترس شیمیدان قابل محاسبه هستند. توصیف‌کننده‌های نظری و یا محاسباتی با استفاده از ساختار ترکیبات شیمیایی و بر اساس معادلات ریاضی ساده و یا پیچیده محاسبه می‌شوند، بنابراین برای همه ترکیبات شیمیایی پس از رسم و بهینه‌سازی ساختار قابل محاسبه هستند و هیچ محدودیت محاسباتی برای آن وجود ندارد. با پیشرفت‌های اخیر در علم شیمی محاسباتی و ایجاد نرم‌افزارهای محاسباتی کارآمد مانند دراگون<sup>۱</sup>، تعداد زیادی توصیف‌کننده مولکولی برای یک مولکول قابل محاسبه است.

## ۱-۵-۵ پیش‌پردازش توصیف‌کننده‌های محاسبه شده

همان‌طور که گفته شد، پس از محاسبه توصیف‌کننده‌ها تعداد زیادی از این توصیف‌کننده‌های مولکولی (متغیرهای مستقل) به‌وجود می‌آید، که همه این‌ها دارای اهمیت و ارتباط ویژه با پاسخ هدف نمی‌باشند. بنابراین قبل از انتخاب مؤثرترین توصیف‌کننده‌ها، عملیات پیش‌پردازش و یا غربالگری<sup>۳</sup> انجام می‌شود. به‌طوری‌که، توصیف‌کننده‌های مولکولی زائد و فاقد اطلاعات مفید از فضای ماتریس متغیرهای مستقل خارج می‌شوند. این فرایند علاوه بر این که سبب کاهش زمان محاسبات در مراحل انتخاب متغیر و مدل‌سازی می‌شود، بلکه صحت پیش‌بینی مدل‌های توسعه یافته را نیز افزایش می‌دهد. به‌منظور انجام فرایند

---

<sup>۱</sup>Dragon

<sup>۲</sup>Pre-processing

<sup>۳</sup>Screening

پیش‌پردازش ابتدا توصیف‌کننده‌هایی با مقادیر ثابت و نسبتاً ثابت با استفاده از نرم‌افزار آماری R [۲۷] و با استفاده از بسته نرم‌افزاری caret [۲۸] حذف می‌شوند. سپس همبستگی بین توصیف‌کننده‌ها نیز با استفاده از دستور corecoeff در نرم‌افزار متلب<sup>۱</sup> مورد بررسی قرار خواهد گرفت. از بین دو توصیف‌کننده با همبستگی بالای ۰/۹، توصیف‌کننده‌ای که بیش‌ترین همبستگی را با متغیر وابسته (فعالیت دارویی یا خاصیت فیزیوشیمیایی) داشت حفظ شده و دیگری حذف خواهد شد.

## ۱-۵-۶ انتخاب توصیف‌کننده‌های مؤثر

روش‌های متفاوتی برای انتخاب مؤثرترین توصیف‌کننده‌ها استفاده می‌شود که می‌توان به روش‌های انتخاب متغیر کلاسیک و روش‌های انقباضی اشاره کرد. روش‌های کلاسیک همچون رگرسیون خطی چندگانه<sup>۲</sup>، دارای مشکلاتی چون بی‌ثباتی، واریانس بالای برآوردگرها، بایاس بالا و تعداد زیاد توصیف‌کننده‌های منتخب و به‌دنبال آن کاهش تفسیرپذیری می‌باشند. از این‌رو روش‌های رگرسیونی انقباضی یا جریمه شده<sup>۳</sup> جدید به هدف حل مشکلات موجود معرفی شدند. روش حداقل قدر مطلق انقباض و عملگر انتخاب‌کننده (LASSO) [۲۹]، لاسوی تطبیقی (ALASSO) [۳۰]، تابع جریمه انحراف قدر مطلق به‌طور هموار بریده شده (SCAD) [۳۱] و حداقل انحراف مطلق - حداقل قدر مطلق انقباض و عملگر انتخاب‌کننده (LAD-LASSO)<sup>۴</sup> [۳۲] به‌عنوان روش‌های انقباضی جدید بسیار کارآمد هستند. در ادامه به معرفی انواع روش‌های انتخاب متغیر پرداخته می‌شود.

---

<sup>۱</sup>MATLAB

<sup>۲</sup>Penalized methods

<sup>۳</sup> Multiple Linear Regression (MLR)

<sup>۴</sup> Shrinkage or penalized regression methods

<sup>۵</sup> Least absolute shrinkage and selection operator (LASSO)

<sup>۶</sup> Adaptive LASSO (ALASSO)

<sup>۷</sup> Smoothly clipped absolute deviation (SCAD)

<sup>۸</sup>Least absolute deviation- Least absolute shrinkage and selection operator

## ۱-۵-۶ انتخاب متغیر به روش رگرسیون خطی چندگانه

رگرسیون، معادله‌ای برای ایجاد یک رابطه بین یک متغیر وابسته از یک طرف و یک یا چند متغیر مستقل از طرف دیگر می‌باشد. اگر تنها یک متغیر مستقل وجود داشته باشد، رگرسیون را ساده و در غیر این صورت، رگرسیون را چندگانه می‌گویند. در بسیاری موارد نمی‌توان تغییرات یک متغیر وابسته را فقط به مقادیر یک متغیر مستقل ارتباط داد. به عبارت دیگر برای پیش‌بینی مقادیر یک متغیر، دانستن مقادیر دو یا چند متغیر دیگر لازم است. زمانی که تعداد متغیرهای مستقل مورد مطالعه زیاد باشد لازم است که تکنیکی برای انتخاب متغیرهای مهم و اثرگذار مورد استفاده قرار گیرد تا بتوان متغیرهای معنی‌دار را از متغیرهای بی‌اهمیت تمییز داد. روش‌های رگرسیون خطی چندگانه متفاوتی برای انتخاب بهترین زیر مجموعه از تمام مجموعه‌ها وجود دارد. از جمله مهم‌ترین آن‌ها می‌توان به انتخاب پیش‌رونده<sup>۱</sup> حذفی پس‌رونده<sup>۲</sup> و روش گام‌به‌گام<sup>۳</sup> اشاره کرد. هر چند این روش‌ها، بهترین زیرمجموعه را از بین تمام زیرمجموعه‌های ممکن انتخاب می‌نمایند، اما زمانی که ابعاد داده‌ها بزرگ باشد تعداد زیر مجموعه‌های منتخب افزایش می‌یابد و با افزایش تعداد متغیرهای مدل، مقدار ضریب تعیین به‌طور کاذب بالا بوده و قدرت پیش‌بینی مدل پایین می‌آید. علاوه بر این، روش‌های انتخاب متغیر کلاسیک دارای مشکل بی‌ثباتی هستند. به‌طوری‌که با تغییری کوچک در داده‌ها، مدل‌های خیلی متفاوتی به وجود می‌آید. با توجه به مشکلات موجود، محققین رگرسیون جریمه شده را معرفی کردند که با اعمال یک پارامتر جریمه<sup>۴</sup> مناسب دو عمل انتخاب متغیر و برآورد ضرایب به‌طور هم‌زمان انجام می‌شود. بنابراین همان‌طور که گفته شد استفاده از روش انتخاب متغیر خطی چندگانه به دلیل داشتن معایبی هم‌چون ناپایداری، بایاس بالا، ناکارآمدی در حضور هم‌خطی، برای انتخاب توصیف‌کننده‌های منتخب توصیه نمی‌شود [۲۹].

---

<sup>1</sup>Forward

<sup>2</sup>Backward

<sup>3</sup>Stepwise

<sup>4</sup>Penalty

مدل رگرسیون خطی زیر را در نظر بگیرید:

$$y_i = x_i' \beta + \epsilon_i \quad i=1, \dots, n$$

که در آن  $x_i = (x_{i1}, \dots, x_{ip})'$  ماتریس توصیف‌کننده‌های مولکولی با ابعاد  $p$  و  $y_i$  متغیر وابسته و  $\beta = (\beta_1, \dots, \beta_p)$  ضرایب رگرسیونی و  $\epsilon_i$  نیز خطای تصادفی با میانه  $0$  می‌باشد. تیبشیرانی<sup>۱</sup> (۱۹۹۶) روشی جدید برای تخمین مدل‌های خطی ارائه کرد و آن را LASSO نامید. این روش به صورت هم‌زمان به برآورد و انتخاب متغیر می‌پردازد. انگیزه اصلی تیبشیرانی در تعریف LASSO، برگرفته از پیشنهاد بریمن<sup>۲</sup> (۱۹۹۳) است.

به پیشنهاد بریمن

$$\sum_{i=1}^n \left\{ y_i - \sum_{j=1}^p c_j \hat{\beta}_j x_{ij} \right\}^2 \quad \text{رابطه ۱-۱}$$

تحت شرایط  $c_j \geq 0$  و  $\sum_{j=1}^k c_j \leq t$  ضریب رگرسیونی حداقل می‌شود که در آن  $t$  یک مقدار ثابت است. اگر  $t=0$  باشد آن‌گاه متغیری در مدل نخواهد بود و اگر  $t = \square$  باشد مدل با مجموعه‌ای از متغیرها ساخته می‌شود.

LASSO را در حالت کلی می‌توان با مینیمم کردن عبارت خطای جریمه شده به صورت زیر به دست آورد:

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n |y_i - x_i' \beta| + \lambda \sum_{i=1}^p |\beta_i| \quad \text{رابطه ۲-۱}$$

که در آن  $\lambda$  یک پارامتر تنظیم‌کننده نامنفی ( $\lambda > 0$ )،  $|\beta| = (|\beta_1|, \dots, |\beta_p|)^T$  می‌باشد. پارامتر  $\lambda$  سطح تنگی و تعداد ضرایب صفر را در برآوردگر LASSO کنترل می‌کند. به عبارت دیگر، هرچه جریمه بزرگ‌تری

<sup>1</sup> Tibshirani

<sup>2</sup> Breiman

به کار برده شود، تعداد بیش‌تری از ضرایب به سمت صفر منقبض می‌شوند. در حقیقت  $\lambda$  و  $t$  رفتار یا رابطه‌ای معکوس با یکدیگر دارند. LASSO با افزایش  $\lambda$  ضرایب را به سمت صفر منقبض می‌کند و زمانی که  $\lambda \rightarrow 0$  ضرایب بیش‌تری در مدل باقی خواهند ماند. در واقع LASSO مجموع توان‌های دوم خطای مدل رگرسیون را تحت این محدودیت که مجموع مقادیر مطلق ضرایب رگرسیونی کم‌تر از یک مقدار ثابت باشند، حداقل می‌سازد. به دلیل اعمال این محدودیت، LASSO تمایل به ایجاد ضرایب دقیقاً برابر با صفر دارد که در نتیجه LASSO می‌تواند به‌عنوان یک روش انتخاب متغیر کارآمد در نظر گرفته شود. اهمیت ذاتی LASSO زمانی بیش‌تر آشکار می‌شود که با داده‌هایی با ابعاد بالا سر و کار داریم. از آنجایی که متغیرهای خیلی کمی ارتباط معنادار با متغیر وابسته دارند، بنابراین استفاده از روشی که بتواند ضرایب متغیرهای بی‌اهمیت را صفر کند مورد نیاز است. تیبشیرانی با معرفی LASSO روشی را معرفی کرد که با آن می‌توان ابعاد داده‌ها را کاهش داد. با توجه به این که تعداد متغیرها کم می‌شود، مدل‌های ساخته شده تفسیرپذیر، ساده و پایدارتر هستند [۲۹].

LASSO پنجره جدیدی را به روی مبحث انتخاب متغیر گشود، به‌طوری که بسیاری از روش‌های انتخاب متغیر که بعد از سال ۱۹۹۶ ارائه شدند به‌نوعی با LASSO در ارتباط هستند و هر یک با هدف اصلاح یکی از نقاط ضعف LASSO ایجاد شدند. در واقع در استفاده از توابع جریمه‌ای، سؤالی که پیش می‌آید، نوع جریمه‌ای است که باید مورد استفاده قرار گیرد. فن و لی<sup>۱</sup> (۲۰۰۱) نشان دادند که یک تابع جریمه خوب باید سه ویژگی مطلوب زیر را داشته باشد.

❖ تنکی: برآوردگر حاصل باید به‌طور خودکار ضریب برآورد شده‌ای که مقدار کوچک دارد را

صفر کند تا متغیرهای مناسب انتخاب شود و به‌این وسیله پیچیدگی مدل کاهش می‌یابد.

---

<sup>1</sup> Fan and Li

<sup>2</sup> Sparsity

❖ نا اریبی: برآوردگر به دست آمده برای ضرایب رگرسیونی با مقادیر بزرگ، تقریباً نا اریب باشد تا اریبی مدل کاهش یابد.

❖ پیوستگی: برآوردگر حاصل پیوسته باشد تا باعث پایداری مدل شود.

با این که LASSO در بسیاری از مسائل عملکرد خوبی را از خود نشان داده است اما دارای محدودیت‌هایی نیز می‌باشد. به عنوان مثال نتیجه جریمه LASSO، برآوردگری با اریبی زیاد می‌باشد، زیرا اندازه انقباضی که به ضرایب رگرسیونی کوچک و بزرگ اختصاص می‌دهد یکسان است و این باعث می‌شود که ضرایب کوچکی که دارای ارتباط مؤثر با متغیر وابسته هستند نیز منقبض و صفر شوند. همچنین LASSO برای یک مدل رگرسیونی خطی با  $p$  متغیر پیشگو و  $n$  مشاهده، حداکثر  $n$  متغیر انتخاب می‌کند. این محدودیت‌ها نشان‌دهنده کارایی پیش‌بینی پایین LASSO در تخمین دقیق ضرایب غیر صفر می‌باشد. منظور از کارایی پیش‌بینی، سازگاری روش در انتخاب متغیر است. به این معنی که یک روش با احتمال نزدیک به یک بتواند ضرایب غیر صفر را به درستی شناسایی کند. از این رو در سال ۲۰۰۱ یک جریمه مناسب با ویژگی‌های نا اریب، تنک و پیوسته به نام SCAD معرفی شد [۳۱] و همچنین لاسوی تطبیقی (ALASSO) با کارایی پیش‌بینی مناسبی در سال ۲۰۰۶ برای رفع ناسازگاری LASSO معرفی شد [۳۰]. روش‌های نامبرده در ادامه به اختصار شرح داده می‌شوند.

### ۱-۵-۶-۳ SCAD

همان‌طور که گفته شد در مقایسه با LASSO، تابع جریمه SCAD، دارای ویژگی‌هایی چون نا اریبی، تنکی و پیوستگی می‌باشد. با توجه به اینکه در روش LASSO ضرایب بزرگ نیز جریمه می‌شوند،

---

<sup>1</sup> Unbiasedness

<sup>2</sup> Continuity



ضرایب کوچک‌تر فرصت حضور در مدل را پیدا می‌کنند و این سبب بیش برآزش<sup>۱</sup> می‌گردد. به این ترتیب LASSO از ویژگی نا اریبی بی بهره است. تابع جریمه SCAD به صورت زیر ارائه شده است:

$$\hat{\beta}^{SCAD} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \|Y - X\beta\|^2 + \lambda \sum_{i=1}^p P(|\beta_i|); a, \lambda \right\} \quad \text{رابطه ۳-۱}$$

که  $P(., a, \lambda)$  تابع جریمه SCAD است و دارای سه آرگومان زیر است:

$$P(t; a, \lambda) = \begin{cases} \lambda t & 0 \leq t \leq \lambda \\ \frac{2a\lambda t - t^2 - \lambda^2}{2(a-1)} & \lambda < t < a\lambda \\ \frac{\lambda^2(a+1)}{2} & t \geq a\lambda \end{cases} \quad \text{رابطه ۴-۱}$$

$a > 2$  و  $\lambda > 0$  پارامترهای جریمه هستند و اندازه انقباض ضرایب را کنترل می‌کنند. با توجه به

رابطه ۳-۱ ملاحظه می‌شود که در فاصله  $|\beta_j| \leq \lambda$  تابع جریمه SCAD بر LASSO منطبق است و پس از آن تا زمانی که  $|\beta_j| < a\lambda$  جریمه SCAD یک تابع درجه دوم است. سپس به ازای  $|\beta_j| \geq a\lambda$  به یک تابع با مقدار ثابت تبدیل می‌شود و وقتی  $\alpha \rightarrow \infty$ ، برآوردگر SCAD با برآوردگر LASSO معادل است. فن و لی با مطالعات متعدد نشان دادند که برای SCAD برابر با  $3/7$  می‌باشد [۳۱]. برنامه SCAD را می‌توان با استفاده از بسته نرم‌افزاری `ncvreg` در برنامه R اجرا نمود [۳۳].

#### ALASSO ۴-۶-۵-۱

به صورت کلی تابع جریمه ALASSO، مانند LASSO است. اما برخلاف LASSO برای ضرایب

مختلف، جریمه‌های متفاوتی را در نظر می‌گیرد.

$$\hat{\beta}_{ALASSO} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \|Y - X\beta\|^2 + \lambda \sum_{i=1}^p \hat{w}_i |\beta_i| \right\} \quad \text{رابطه ۵-۱}$$

<sup>1</sup> Overfitting

$\hat{w}_i$  وزن بردار می‌باشد. با توجه به معادله فوق، اگر وزن‌ها به‌درستی انتخاب شوند این نوع از

LASSO دارای کارایی پیش‌بینی بالاتری خواهد بود. وزن‌ها نیز به‌صورت زیر تعریف می‌شوند:

$$\hat{w}_i = \frac{1}{|\hat{\beta}_i|^\gamma} \quad \gamma > 0 \quad \text{رابطه ۶-۱}$$

که در آن  $\beta_i$  می‌تواند برآوردگر رگرسیون کم‌ترین توان‌های دوم، LASSO و یا Ridge باشد. زو

ALASSO (۲۰۰۶) را به هدف رفع ناسازگاری مدل LASSO معرفی کرد. ALASSO با توجه به اعمال

جریمه‌های متفاوت روی ضرایب با بزرگی متفاوت، سبب بهبود کارایی پیش‌بینی مدل شد. ALASSO را

می‌توان با استفاده از بسته نرم‌افزاری parcore در برنامه R اجرا نمود [۳۴].

#### ۱-۵-۶-۵ LAD-LASSO

در مدل رگرسیون خطی

$$y_i = x_i' \beta + \epsilon_i \quad i=1, \dots, n \quad \text{رابطه ۷-۱}$$

برای برآورد پارامترها از حداقل سازی عبارت OLS که برابر است با  $\sum_{i=1}^n (y_i - x_i' \beta)^2$  استفاده می‌شود.

اما عملکرد برآوردهای رگرسیونی OLS، با وجود داده‌های دور افتاده محدود می‌شود. برای مقابله با

مشاهده‌های دور افتاده روش‌های رگرسیونی استواری چون برآوردگر حداقل قدر مطلق انحراف<sup>۱</sup> (LAD)

پیشنهاد می‌شود.

$$\hat{\beta}_{LAD} = \operatorname{argmin} \sum_{i=1}^n |y_i - x_i' \beta| \quad \text{رابطه ۸-۱}$$

در سال ۱۹۹۶ تیبشیرانی برای صفر کردن ضرایب بی‌اهمیت LASSO را معرفی نمود که با نگاهی

به رابطه ۲-۱ مشاهده می‌شود که با اضافه شدن عبارت جریمه به معادله حداقل مربعات معمولی ضرایبی با

<sup>۱</sup>Least Absolute Deviations

مقدار ناچیز صفر می‌شوند. هر چه مقدار  $\lambda$  بزرگ‌تر باشد میزان انقباض بیش‌تر و ضرایب رگرسیون کم‌تری باقی می‌ماند و در غیر این‌صورت با کوچک شدن مقدار  $\lambda$  تعداد ضرایب غیر صفر بیش‌تری باقی می‌ماند. با توجه به اینکه LASSO برای همه ضرایب از یک مقدار  $\lambda$  واحد برای انقباض استفاده می‌کند پارامتر تنظیم یکسانی برای همه پارامترها اعمال می‌شود و مدل از بایاس و اریبی رنج می‌برد. با اضافه شدن تابع جریمه نرم L1 به برآوردگر LAD می‌توان یک مدل تنک به‌دست آورد که علاوه بر مواجهه با بایاس، به اریبی در برابر داده‌های دور افتاده، نیز حساسیت ندارد. این روش در مقایسه با LAD می‌تواند عمل انتخاب متغیر و برآورد پارامتر را به‌صورت هم‌زمان انجام دهد و در مقایسه با LASSO در مقابل مشاهده‌های دورافتاده استوار است.

$$\hat{\beta}_{\text{LAD-LASSO}} = \operatorname{argmin} \sum_{i=1}^n |y_i - x_i' \beta| + \lambda_j \sum_{j=1}^p |\beta_j| \quad \text{رابطه ۹-۱}$$

برآوردگر LAD-LASSO را می‌توان به‌عنوان یک برآوردگر بیزی در نظر گرفت، به‌طوری که هر ضریب رگرسیونی  $\beta_j$  دارای پارامتر مقیاسی  $\lambda_j n$  است و بر این اساس مقدار  $\hat{\lambda}_j = \frac{1}{n|\beta_j|}$  است که با استفاده از برآورد LAD معمولی برآورد می‌شود، که با استفاده از تکنیک تخمین LAD معمولی تخمین زده می‌شود. LAD-LASSO را می‌توان با استفاده از بسته نرم‌افزاری MTE در برنامه R بدون برنامه ریزی محاسباتی پیچیده به‌دست آورد [۳۵].

## ۱-۵-۶ ارزیابی توصیف‌کننده‌های منتخب

توصیف‌کننده‌های منتخب روش انتخاب متغیر مورد نظر باید از نظر وجود همبستگی و هم‌خطی مورد ارزیابی قرار گیرند. به‌عبارت دیگر توصیف‌کننده‌های مولکولی منتخب باید فاقد همبستگی بالای ۰/۹ (معنی‌دار) و فاقد هم‌خطی باشند. پدیده همبستگی با استفاده از محاسبه مربع ضریب همبستگی به‌دست

می‌آید. با رسم نمودار نقشه رنگی<sup>۱</sup> نتایج با کیفیت و وضوح بیشتری قابل نمایش خواهد بود. مقادیر همبستگی نزدیک به ۱ و ۱- نشاندهنده همبستگی بالا بین دو توصیف کننده می‌باشد. پدیده هم‌خطی با استفاده از محاسبه پارامتر افزایش عامل واریانس<sup>۲</sup> (VIF) مورد بررسی قرار گرفت. VIF مربوط به هر توصیف کننده منتخب با استفاده از رابطه ۱-۱۰ محاسبه شد.

$$VIF = \frac{1}{1 - R_i^2} \quad i=1,2,3,\dots,p \quad \text{رابطه ۱-۱۰}$$

p در رابطه ۱-۱۰ تعداد پارامترهای منتخب روش انتخاب متغیر و  $R_i^2$  نیز مجذور همبستگی چندگانه است که از رگرسیون یک متغیر بر سایر متغیرها به دست می‌آید. قرارگیری مقدار VIF در محدوده ۱-۵ دلالت بر عدم وجود پدیده هم‌خطی نگران کننده بین توصیف کننده‌های منتخب می‌باشد [۳۶-۳۸].

## ۱-۵-۷ مدل سازی

به منظور ایجاد ارتباط بین متغیرهای مستقل و وابسته، مدل‌های ریاضی، ساخته می‌شود. با استفاده از مدل‌های توسعه یافته به راحتی می‌توان پاسخ هدف مربوط به ترکیبات جدید را تخمین زد. با پیشرفت علم و تکنولوژی، روش‌های مدل‌سازی متفاوتی گسترش یافته‌اند که از آن جمله می‌توان به رگرسیون خطی چندگانه (MLR)، رگرسیون اجزای اصلی (PCR)، کم‌ترین توان‌های دوم جزئی (PLS)؛ شبکه عصبی مصنوعی (ANN) اشاره کرد. در این پژوهش از روش شبکه عصبی مصنوعی استفاده شده است.

## ۱-۵-۷-۱ شبکه عصبی مصنوعی

شبکه عصبی یک برنامه نرم‌افزاری است که می‌تواند همانند مغز انسان عمل نماید. در واقع یک شبکه عصبی مصنوعی ایده‌ای است برای پردازش اطلاعات که از سیستم عصبی زیستی الهام گرفته شده

<sup>۱</sup>Heatmap

<sup>۲</sup>Variance inflation factor

<sup>۳</sup>Principal component regression

<sup>۴</sup>Partial least square

<sup>۵</sup>Artificial neural network

و مانند مغز انسان به پردازش اطلاعات می پردازد. این سیستم‌ها از تعداد زیادی عنصر پردازش به نام نورون تشکیل شده‌اند که برای حل یک مسئله به صورت هماهنگ با هم عمل می کنند. شبکه‌های عصبی مصنوعی نظیر مغز انسان‌ها، با مثال یاد می گیرند و با پردازش روی داده‌های تجربی، دانش یا قانون نهفته در ورای داده‌ها را به ساختار شبکه منتقل می کنند. به همین خاطر به این سیستم‌ها، هوشمند گفته می شود، زیرا شبکه‌ها بر اساس محاسبات روی داده‌های عددی یا مثال‌ها، قوانین کلی را فرا می گیرند. امروزه از شبکه‌های عصبی به عنوان یک ابزار کارآمد در زمینه‌های مختلف علمی از جمله صنایع الکترونیک، پزشکی، اکتشاف نفت و گاز، رباتیک، شیمی و داروسازی استفاده می شود. در این پروژه از روش شبکه عصبی به عنوان یک مدل قدرتمند غیرخطی برای ایجاد ارتباط بین متغیرهای مستقل و وابسته استفاده شده است.

### -آموزش شبکه‌های پیشخور با تکنیک پس انتشار<sup>۱</sup>خطا

تکنیک پس انتشار خطا یک روش متداول آموزش با ناظر برای شبکه‌های پیشخور 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 پس انتشار طبق مراحل زیر انجام می شود [۳۹].

۱- انتشار ورودی‌ها از نورون‌های ورودی به سمت نورون‌های خروجی

۲- اختصاص ماتریس وزن‌های تصادفی به هریک از اتصالات

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

۴- پس انتشار خطا از نورون‌های خروجی به سمت نورون‌های ورودی و اصلاح وزن‌ها

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

---

<sup>۱</sup>Back propagation

مراحل فوق تا رسیدن به حداکثر تکرار<sup>۱</sup> (دور آموزش) مجاز انجام می‌شود و یا این که مقدار تابع کارآیی از مقداری که تعیین شده کم‌تر باشد. شبکه عصبی انتخاب شده در این پروژه یک شبکه پیشخور با الگوریتم آموزشی پس انتشار خطا می‌باشد [۳۹].

در این مطالعه، یک ANN پیشخور با الگوریتم پس انتشار خطا برای مدل‌سازی QSAR/QSPR در نظر گرفته شده است. در اکثر تحقیقات در زمینه شیمی محاسباتی، ساخت مدل ANN با یک لایه کافی است [۴۰]. بنابراین، مدل‌های ANN توسعه یافته، دارای سه لایه شامل یک لایه ورودی، یک لایه پنهان و یک لایه خروجی است که برای بهینه‌سازی پارامترهای ANN استفاده می‌شود. برای به دست آوردن بهترین مدل شبکه عصبی مصنوعی، تمامی پارامترهای مؤثر بر عملکرد پیش‌بینی مدل شبکه عصبی مصنوعی مانند تعداد ورودی‌ها، گره‌ها در لایه پنهان، دوره‌های آموزشی و توابع آموزش و انتقال بهینه‌سازی می‌شوند. به این منظور، چهار مدل شبکه عصبی مصنوعی مختلف با استفاده از دو الگوریتم آموزشی متفاوت (لونبرگ – مارکواریت (LM)<sup>۲</sup> و تنظیم بیزین (BR)<sup>۳</sup> و دو تابع انتقال متفاوت برای لایه پنهان (لگاریتم سیگموئیدی و تانژانت هایپربولیک سیگموئیدی) برای همه مطالعات QSAR/QSPR طراحی می‌شود. در تمام مدل‌های شبکه عصبی مصنوعی، تابع خطی به عنوان تابع انتقال خروجی استفاده می‌شود و پارامترهای دیگری مانند تعداد ورودی، تعداد گره‌ها و تعداد دوره‌های آموزشی به طور هم‌زمان بهینه می‌شوند. در نهایت از بین همه ساختارهای شبکه عصبی تعیین شده، مدل شبکه عصبی بهینه با استفاده از کم‌ترین مقدار MSE/RMSE مربوط به مجموعه ارزیابی انتخاب می‌شود و از مدل شبکه عصبی بهینه برای پیش‌بینی مقادیر پاسخ ترکیبات مجموعه آزمون استفاده خواهد شد.

---

<sup>۱</sup>Epoch

<sup>۲</sup>Levenberg –Marquardt

<sup>۳</sup>Bayesian regularization

## ۱-۵-۷-۲ چیدمان توصیف‌کننده‌های منتخب به‌عنوان ورودی مدل شبکه عصبی مصنوعی

توصیف‌کننده‌های مهم و دارای ضرایب غیر صفر انتخاب شده به‌وسیله روش‌های رگرسیون انقباضی بر اساس ترتیب ورودی ماتریس توصیف‌کننده‌ها انتخاب می‌شوند. از آنجایی که برای مدل‌سازی شبکه عصبی مصنوعی و بهینه‌سازی پارامترهای شبکه عصبی ترتیب ورودی توصیف‌کننده‌های منتخب حائز اهمیت است، بنابراین باید برای این مسئله، تدبیری اندیشه شود تا بهترین زیر مجموعه به مدل معرفی شود. با توجه به تعداد زیاد این زیرمجموعه‌ها که از حاصل ضرب تعداد حالات توصیف‌کننده‌ها (دوتایی، سه‌تایی،  $n-1$  تایی که  $n$  برابر با تعداد توصیف‌کننده‌های منتخب است) در تعداد گره و تعداد دور آموزشی به‌دست می‌آید، طراحی و استفاده از تمام تولیدات تصادفی به‌عنوان ورودی ANN عملاً غیرممکن است. از این‌رو در این مطالعه از دو روش برای چیدمان توصیف‌کننده‌های منتخب در مدل ANN استفاده می‌شود.

### - چیدمان توصیف‌کننده‌ها بر اساس بزرگی ضرایب استاندارد شده

در روش اول ضرایب استاندارد نشده توصیف‌کننده‌ها با استفاده از اجرای روش‌های رگرسیونی انقباضی به‌دست می‌آید و در نهایت، ضرایب استاندارد نشده در نسبت بین انحراف معیار متغیر مستقل مربوطه و انحراف معیار متغیر وابسته ضرب می‌شود تا ضرایب رگرسیونی از بزرگی مقادیر توصیف‌کننده‌ها مستقل شود. بنابراین ضرایب استاندارد شده با استفاده از رابطه ۱-۱۱ محاسبه می‌شود [۴۱]:

$$\beta_{\text{شده استاندارد}} = \beta_{\text{نشده استاندارد}} \times \frac{S_{\text{مستقل متغیر}}}{S_{\text{وابسته متغیر}}} \quad \text{رابطه ۱-۱۱}$$

در رابطه ۱-۱۱ صورت کسر نشان‌دهنده انحراف استاندارد محاسبه شده برای مقادیر هر توصیف‌کننده به‌عنوان متغیر مستقل و مخرج کسر به‌معنی انحراف استاندارد مربوط به متغیر وابسته می‌باشد. پس از محاسبه ضرایب استاندارد شده، توصیف‌کننده‌ها بر اساس بزرگی قدر مطلق ضرایب استاندارد شده چیده و به‌عنوان ورودی شبکه عصبی مصنوعی تعریف و بهینه‌سازی می‌شوند.

## - پچیدمان توصیف کننده‌ها بر اساس اهمیت در شبکه عصبی مصنوعی

روش دیگر برای ایجاد چیدمان منطقی، استفاده از روش شبکه عصبی است. به طوری که ابتدا مدل شبکه عصبی با تعداد کل توصیف کننده‌های منتخب آموزش داده می‌شود و پارامترهای دیگر شبکه عصبی از جمله گره، دور آموزشی، تابع انتقال و تابع آموزش بهینه به دست می‌آید. سپس برای تخمین اهمیت توصیف کننده  $i$  ام، همه توصیف کننده‌های منتخب روش انتخاب متغیر انقباضی تصادفی می‌شوند. مدل شبکه عصبی بهینه هر بار با استفاده ماتریس ورودی دست کاری شده در حضور توصیف کننده  $i$  با مقادیر تصادفی، آموزش داده می‌شود و مقادیر پاسخ مجموعه ارزیابی پیش‌بینی می‌شود. مقدار  $MSE/RMSE$  با استفاده از پاسخ‌های پیش‌بینی شده و واقعی مجموعه ارزیابی محاسبه می‌شود. این فرآیند تا زمانی که مقادیر هر کدام از توصیف کننده‌ها یکبار با مقادیر تصادفی جایگزین و مقدار خطای ارزیابی محاسبه شود، تکرار می‌شود. در نتیجه، به تعداد توصیف کننده‌های منتخب، مقدار  $MSE/RMSE$  به دست خواهد آمد. بالاترین  $MSE/RMSE$  به این معنی است که مدل ANN در غیاب آن توصیف کننده با مقادیر واقعی، خطای بیش‌تری را متحمل می‌شود و چنین توصیف کننده‌ای بیش‌ترین اهمیت را در ساخت مدل دارد. بنابراین، توصیف کننده‌ها بر اساس اهمیت محاسبه شده (مقدار  $MSE/RMSE$ ) مرتب می‌شوند. در نهایت پارامترهای مدل شبکه عصبی با استفاده از توصیف کننده‌های چیده شده بر اساس اهمیت در شبکه عصبی مصنوعی بهینه می‌شوند. برای درک بهتر مسئله، لازم به ذکر است که زیرمجموعه اول شامل توصیف کننده‌های مهم اول و دوم (دو توصیف کننده با بیش‌ترین  $MSE/RMSE$  مجموعه ارزیابی مدل شبکه عصبی) است و زیرمجموعه‌های بعدی با افزودن توصیف کننده‌های دیگر بر اساس ترتیب اهمیت‌شان به زیرمجموعه اول ایجاد می‌شوند. بنابراین، در بهینه‌سازی شرایط شبکه عصبی مصنوعی، تعداد ورودی‌ها با استفاده از زیرمجموعه‌هایی شامل ۲ تا تعداد کل توصیف کننده‌های منتخب به‌عنوان ورودی شبکه عصبی مصنوعی، تعریف شده و به‌طور هم‌زمان همراه با سایر پارامترهای شبکه عصبی بهینه می‌شوند. در این رساله مدل با



استفاده از هر دو روش تعیین اهمیت توصیف کننده‌ها، توسعه یافته‌اند و نتیجه مربوط به بهترین روش به‌عنوان اساس تعیین اهمیت برتر گزارش شده است و سایر مراحل QSAR/QSPR ادامه یافته است.

## ۱-۵-۸ ارزیابی مدل

هدف اصلی هر مدل‌سازی QSAR/QSPR قابل اعتماد و دقیق، ایجاد یک مدل قوی با ظرفیت بالای پیش‌بینی فعالیت ترکیبات جدید است. از این‌رو ارزیابی مدل توسعه یافته در مطالعات QSAR/QSPR از مراحل اصلی ساخت مدل به‌شمار می‌رود. ارزیابی مدل از طریق شاخص‌های کمی و آنالیزهای آماری متفاوتی انجام پذیر است. از دسته روش‌های ارزیابی مدل می‌توان به محاسبه پارامترهای آماری برای مجموعه ارزیابی خارجی (آزمون) و مجموعه پیش‌بینی شده به‌وسیله تکنیک رد مرحله‌ای تک تک<sup>۱</sup> (LOO)، نمودار باقی‌مانده‌های استاندارد شده، نمایش دامنه کاربرد<sup>۲</sup> مدل و آزمون Y-تصادفی اشاره کرد.

### ۱-۵-۸-۱ بررسی مدل با استفاده از نتایج پیش‌بینی شده برای مجموعه آزمون

برای ارزیابی مدل توسعه یافته، از پیش‌بینی داده‌های مجموعه آزمون که در هیچ یک از مراحل انتخاب متغیر و مدل‌سازی حضور نداشته‌اند، استفاده می‌شود. به‌طوری‌که ترکیبات مجموعه ارزیابی خارجی (آزمون) با استفاده از مدل‌های شبکه عصبی بهینه، پیش‌بینی می‌شوند. مقادیر خطای حاصل از پیش‌بینی و نتایج پارامترهای آماری قابل قبول، صحت و دقت مدل ساخته شده شبکه عصبی جفت شده با روش‌های انتخاب متغیر انقباضی را اثبات می‌کند.

---

<sup>۱</sup>Leave-one-out

<sup>۲</sup>Applicability Domain

## ۱-۵-۸-۲ ارزیابی مدل با استفاده از تکنیک رد مرحله‌ای تک تک

تکنیک LOO با به‌کارگیری کل مجموعه داده‌ها برای بررسی استواری<sup>۱</sup> مدل به‌کار گرفته می‌شود، به این معنی که پایداری مدل را با تغییر در داده‌های مدل اثبات می‌کند. بنابراین شبکه عصبی بهینه برای پیش‌بینی همه داده‌ها، مورد استفاده قرار می‌گیرد. با این تفاوت که در این تکنیک هر داده مورد مطالعه یک‌بار به‌عنوان داده آزمون خارج شده و مدل با مابقی داده‌ها آموزش داده می‌شود و پیش‌بینی پاسخ مربوط به داده خارج شده با استفاده از مدل بهینه به‌دست می‌آید. این عملیات به تعداد کل داده‌ها تکرار شده تا پاسخ همه داده‌ها یک‌بار پیش‌بینی شود. با محاسبه پارامترهای آماری برای نتایج حاصل از تکنیک رد مرحله‌ای تک استواری و استحکام مدل‌های QSAR و QSPR توسعه یافته تأیید می‌شود.

## ۱-۵-۸-۳ نمودار باقی‌مانده‌های استاندارد شده

علاوه بر روش‌های ارزیابی یاد شده، نمودار باقی‌مانده استاندارد شده نیز، برای بررسی نتایج پیش‌بینی شده با استفاده از مدل‌های QSAR/QSPR بهینه، مورد استفاده قرار می‌گیرد. بنابراین مقادیر باقی‌مانده‌های استاندارد شده ( $r_i$ ) با استفاده از رابطه ۱-۱۲، برای نتایج پیش‌بینی شده مجموعه آزمون و تکنیک LOO به‌طور مجزا محاسبه شده و نمودار باقی‌مانده‌های استاندارد شده، از رسم  $r_i$  بر حسب مقادیر واقعی پاسخ به‌دست می‌آید.

$$r_i = \frac{e_i}{s_{e_i}} = \frac{(y_i - \hat{y}_i)}{s_{e_i}} \quad \text{رابطه ۱-۱۲}$$

که در آن  $e_i$  تفاوت بین پاسخ‌های واقعی و پیش‌بینی شده برای هر مشاهده است و  $i = 1, \dots, n$  و  $s_{e_i}$  انحراف استاندارد مقادیر باقی‌مانده است. توزیع یکنواخت داده‌ها حول محور صفر نشان‌دهنده عدم وجود خطای سیستماتیک در مدل‌های توسعه یافته است و برازش مناسب مدل‌های توسعه یافته QSAR/QSPR را ثابت می‌کند.

<sup>۱</sup>Robustness

## ۱-۵-۸-۴ پارامترهای آماری

پارامترهای آماری دسته‌ای از شاخص‌های کمی آماری هستند که به وسیله آن‌ها صحت نتایج ارائه شده توسط مدل اندازه‌گیری می‌شود. پارامترهای آماری محاسبه شده در این رساله در جدول ۱-۱ آورده شده است. پارامترهای آماری از نوع خطا دارای محدوده پذیرش خاصی نیستند، بلکه هر چقدر پارامترهای خطای مدل‌های توسعه یافته کم‌تر باشد، آن مدل از برتری قابل توجهی نسبت به سایر مدل‌ها برخوردار است. از جمله پارامترهای ارزیابی دیگر می‌توان به ضرایب تعیین  $R^2$  و  $Q_{L00}^2$  اشاره کرد. مقادیر قابل قبول ضرایب رگرسیونی  $R^2$  و  $Q_{L00}^2$ ، نشان‌دهنده قدرت پیش‌بینی مدل هستند. به طوری که هر چه این پارامترها به مقدار ۱ نزدیک‌تر باشند، نشان‌دهنده این است که مقادیر پیش‌بینی شده با خطای کمی نزدیک به مقادیر واقعی می‌باشند. مقدار  $R^2$  بزرگ‌تر از ۰/۶ و نزدیک به ۱ از مشخصه‌های لازم برای یک مدل QSAR/QSPR توسعه یافته می‌باشد. از نظر فرمول محاسباتی با پارامتر  $R^2$  برابری دارد و تنها تفاوت آن در نوع پیش‌بینی مقادیر پاسخ است که با استفاده از روش ارزیابی تقاطعی<sup>۲</sup> به روش رد مرحله‌ای تک تک (بخش ۱-۵-۸-۲) پیش‌بینی می‌شود. مقدار  $Q_{L00}^2$  بزرگ‌تر از ۰/۵ و نزدیک‌تر به ۱، برای تأیید اعتمادپذیری و قدرت پیش‌بینی مناسب مدل، ضروری است.

### -محاسبه پارامترهای آماری تروپشا<sup>۳</sup>

محققین در یافته‌اند که پارامترهای خطا و مقادیر ضرایب رگرسیونی قابل قبول، به تنهایی نمی‌توانند برتری و قدرت پیش‌بینی مدل‌های QSAR/QSPR توسعه یافته را اثبات نمایند [۴۲، ۴۳]. از این رو تروپشا، پارامترهای ارزیابی دیگری که وابسته به ضریب تعیین ( $R^2$ ) مدل برتر است را پیشنهاد کرد. همان‌طور که

<sup>1</sup>Determination Coefficients

<sup>2</sup>Cross validation

<sup>3</sup>Tropsha

گفته شد،  $R^2$  به عنوان یک پارامتر ارزیابی پر کاربرد، از رسم پاسخ‌های پیش‌بینی شده بر حسب پاسخ‌های تجربی به دست می‌آید. پارامتر  $R_0^2$  ضریب تعیین حاصل از رسم پاسخ‌های پیش‌بینی شده بر حسب پاسخ‌های تجربی با عرض از مبدأ صفر است.  $R_0'^2$  ضریب تعیین حاصل از رسم پاسخ‌های تجربی بر حسب پاسخ‌های پیش‌بینی شده با عرض از مبدأ صفر است. از طرفی،  $k$  شیب معادله رگرسیون برای نمودار پاسخ‌های پیش‌بینی شده بر حسب پاسخ‌های تجربی با عرض از مبدأ صفر و  $k'$  شیب معادله رگرسیون برای نمودار پاسخ‌های تجربی بر حسب مقادیر پاسخ‌های پیش‌بینی شده با عرض از مبدأ صفر است و وجود شیب معادله رگرسیون بین دو حد  $0/85$  و  $1/15$  حاکی از برتری مدل توسعه یافته QSAR/QSPR است.  $R_{0,R}^2$  و  $R_{0,R}'^2$  نشان‌دهنده میزان نزدیکی پارامترهای  $R^2, R_0^2$  و  $R'^2, R_0'^2$  به هم هستند. نزدیکی این پارامترها به ضریب تعیین مدل و تفاضل کم‌تر از  $0/1$ ، اعتماد پذیری و تعمیم‌پذیری بالای مدل را اثبات می‌کند. فرمول محاسباتی همه پارامترهای یاد شده، در جدول ۱-۱ آورده شده است.

### – محاسبه پارامترهای آماری روی<sup>۱</sup>

محقق دیگری به نام روی، سال‌ها بعد پارامترهای آماری تروپشا را گسترش داد و برخی پارامترهای دیگر را از روی پارامترهای تروپشا معرفی نمود [۴۴]. پارامترهای  $R_m^2$  و  $R_m'^2$  پارامترهای آماری دیگری هستند که اختلاف بین  $R^2$  و  $R_0^2$  را در مدل نشان می‌دهند. اگر مدل توسعه یافته دارای  $R_m^2$  بیش‌تر از  $0/5$  باشد، مدل از اعتماد پذیری مناسب و قدرت پیش‌بینی خوبی برخوردار است. فرمول محاسباتی این پارامتر در جدول ۱-۱ آورده شده است.  $R-R$  پارامتر حاصل از تفاضل دو پارامتر  $R_m^2$  و  $R_m'^2$  است. مقدار کوچک‌تر از  $0/3$  این پارامتر قابل قبول است و تأیید دیگری بر مدل QSAR/QSPR توسعه یافته است.

<sup>1</sup>Roy

جدول ۱-۱ پارامترهای آماری

ردیف	پارامتر آماری	فرمول محاسباتی	محدوده قابل قبول
۱	PRESS	$\sum (y_i - \hat{y}_i)^2$	-
۲	SEP	$\sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y})^2}{n}}$	-
۳	MAE	$\frac{\sum_{i=1}^n  y_i - \hat{y} }{n}$	-
۴	REP(%)	$\frac{100}{\bar{y}} \times \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y})^2}{n}}$	-
۵	MSE	$\frac{\sum_{i=1}^n (y_i - \hat{y})^2}{n}$	-
۶	MRE	$\frac{\sum_{i=1}^n \left  \frac{y_i - \hat{y}}{y_i} \right }{n} \times 100$	-
۷	R <sup>2</sup>	$1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$	R <sup>2</sup> > ۰.۱۶
۸	Q <sup>2</sup> <sub>Loo</sub>	R <sup>2</sup> مربوط به داده‌های پیش‌بینی شده بر اساس ارزیابی تقاطعی به تکنیک LOO	Q <sup>2</sup> <sub>Loo</sub> > ۰.۱۵
۹	R <sup>2</sup> <sub>0</sub>	ضریب تعیین داده‌های پیش‌بینی شده بر حسب تجربی با عرض از مبدأ صفر	نزدیک به R <sup>2</sup>
۱۰	R <sup>2</sup> <sub>0.R</sub> نسبی	$R_{0.R}^2 = \frac{(R^2 - R_0^2)}{R^2}$	< ۰.۱
۱۱	R <sup>2</sup> <sub>m</sub>	$R^2 \times [1 - (R^2 - R_0^2)^{\frac{1}{2}}]$	> ۰.۱۵
۱۲	R' <sup>2</sup> <sub>0</sub>	ضریب تعیین داده‌های تجربی بر حسب پیش‌بینی شده با عرض از مبدأ صفر	R <sup>2</sup> به
۱۳	R' <sup>2</sup> <sub>0.R</sub> نسبی	$R_{0.R}'^2 = \frac{(R^2 - R_0'^2)}{R^2}$	< ۰.۱
۱۴	R' <sup>2</sup> <sub>m</sub>	$R^2 \times [1 - ((R^2 - R_0'^2))^{\frac{1}{2}}]$	> ۰.۱۵
۱۵	R-R	$ R_m^2 - R_m'^2 $	< ۰.۳
۱۶	k	شیب نمودار داده‌های پیش‌بینی شده بر حسب تجربی با عرض از مبدأ صفر	۰.۱۸۵ ≤ k ≤ ۱/۱۵
۱۷	k'	شیب نمودار داده‌های تجربی بر حسب پیش‌بینی شده با عرض از مبدأ صفر	۰.۱۸۵ ≤ k' ≤ ۱/۱۵

## ۱-۵-۸-۵ دامنه کاربرد مدل

یکی دیگر از روش‌های ارزیابی استحکام و اعتمادپذیری مدل‌های QSAR/QSPR توسعه یافته، آنالیز دامنه کاربرد است. دامنه کاربرد، یک فضای شیمیایی تئوری است که با استفاده از توصیف‌کننده‌های مولکولی مجموعه آموزش و پاسخ تجربی و پیش‌بینی شده مربوطه ایجاد می‌شود. بنابراین اگر داده شیمیایی جدید نیز در این فضای شیمیایی تئوری قرار گیرد، نشان‌دهنده این است که مدل توانسته در مواجهه با داده‌ای که تاکنون در مدل حضور نداشته، آن را به خوبی پیش‌بینی نماید. دامنه کاربرد با استفاده از محاسبه پارامتر Leverage (H) مجموعه داده‌های آموزش به دست می‌آید. مقادیر H هر ترکیب با استفاده از انعکاس پاسخ واقعی بر پاسخ پیش‌بینی شده محاسبه می‌شود. مقادیر H با توجه به رابطه ۱-۱۳ قابل محاسبه است:

$$H = x_i (X^T X)^{-1} x_i^T \quad \text{رابطه ۱-۱۳}$$

که در اینجا X ماتریس توصیف‌کننده‌های مولکولی داده‌های مجموعه آموزش می‌باشد و  $x_i$  بردار توصیف‌کننده‌های مربوط به هر داده است. T نیز نشان‌دهنده واریانس ماتریس است. برای نمایش فضای شیمیایی دامنه کاربرد از نمودار ویلیام استفاده می‌شود. به طوری که مقادیر باقی‌مانده‌ها با استفاده از مقادیر پیش‌بینی شده توسط مدل برتر و مقادیر واقعی محاسبه می‌شود. در نهایت مقادیر باقی‌مانده‌ها با استفاده از رابطه زیر استاندارد می‌گردد.

$$r_i = \frac{(y_i - \hat{y}_i) - \bar{y}}{s_{e_i}} \quad \text{رابطه ۱-۱۴}$$

که در آن  $y_i$  و  $\hat{y}_i$  به ترتیب نشان‌دهنده پاسخ‌های واقعی و پیش‌بینی شده برای هر مشاهده است و  $s_{e_i}$  انحراف استاندارد مقادیر باقی‌مانده و  $\bar{y}$  میانگین مقادیر واقعی پاسخ است. بنابراین، نمودار ویلیام از رسم مقادیر باقی‌مانده‌های استاندارد شده در مقابل مقادیر H رسم می‌شود. برای آنالیز دامنه کاربرد، داده‌های شیمیایی باید در دو محدوده اطمینان قابل قبول نمودار ویلیام قرار گیرند. اولین شرط قابل قبول،

---

<sup>1</sup>William's plot

قرارگیری و عدم تجاوز مقدار باقی مانده‌های استاندارد شده داده‌های شیمیایی در محدوده ۳ برابری بزرگ‌تر/ کوچک‌تری از انحراف استاندارد می‌باشد. علاوه بر این شرط، مقادیر محاسبه شده H نیز نباید از مقدار حد آستانه یا حد هشدار  $h^*$  برابر با  $3p/n$  بزرگ‌تر باشند. p در این معادله برابر با تعداد توصیف‌کننده‌های مدل به علاوه یک و n تعداد داده‌های مجموعه آموزش می‌باشد. بنابراین، قرارگیری داده‌های H، در محدوده‌های قابل قبول، استحکام و قابل اعتماد بودن مدل‌های QSAR/QSPR توسعه یافته را اثبات می‌کند [۴۵، ۴۶].

### ۱-۵-۸-۶ آزمون Y-تصادفی

آزمون Y-تصادفی برای بررسی استحکام مدل و عدم وجود ارتباط تصادفی بین متغیرهای مستقل و متغیر وابسته به کار گرفته می‌شود. از این‌رو مقادیر پاسخ هدف در محدوده تغییرات پاسخ، ۱۰۰۰ بار تصادفی می‌شود. مدل شبکه عصبی بهینه با پاسخ‌های تصادفی توسعه داده می‌شود و پس از ۱۰۰۰ بار اجرا، پاسخ‌های مربوط به هر مجموعه با استفاده از مدل بهینه پیش‌بینی می‌شود. یک مدل QSAR/QSPR مطمئن باید دارای ضرایب تعیین بسیار کوچک‌تر از مقدار قابل قبول  $0/6$  در مدل‌های تصادفی باشد. مقادیر بسیار کوچک‌تر  $R^2$  نشان‌دهنده این است که ارتباط بین توصیف‌کننده‌ها و مقادیر واقعی پاسخ، شانسی و تصادفی نبوده و یک رابطه منطقی ریاضی بین آن‌ها برقرار است و مدل توسعه یافته بهینه از استحکام و اعتمادپذیری مناسبی برخوردار است.

### ۱-۵-۸-۷ بررسی مشارکت توصیف‌کننده‌های منتخب در شبکه عصبی

با توجه به مدل غیرخطی توسعه یافته، سهم مشارکت توصیف‌کننده‌های منتخب در مدل نهایی با استفاده از مدل بهینه شبکه عصبی پیشنهادی مورد بررسی قرار می‌گیرد. برای محاسبه درصد مشارکت هر توصیف‌کننده موجود در مدل بهینه، مقادیر هر توصیف‌کننده در محدوده تغییرات خاص خود تصادفی می‌شود. مدل‌های بهینه هر بار در حضور یک توصیف‌کننده با مقادیر تصادفی و سایر توصیف‌کننده‌ها با مقادیر واقعی خود توسعه داده می‌شوند. مدل‌های توسعه یافته برای پیش‌بینی مقادیر پاسخ مجموعه ارزیابی

مورد استفاده قرار خواهند گرفت. مقدار پارامتر میانگین قدر مطلق خطا (MAE) مجموعه ارزیابی در حضور توصیف کننده با مقادیر تصادفی، محاسبه می‌شود. این فرآیند برای همه توصیف کننده‌ها تکرار شده و به تعداد توصیف کننده‌های مدل MAE محاسبه می‌شود. درصد مشارکت هر توصیف کننده ( $C_i$  %)، با استفاده از رابطه ۱-۱۵ محاسبه می‌شود.

$$C_i = \left( \frac{MAE_i}{\sum MAE_i} \right) \times 100 \quad \text{رابطه ۱-۱۵}$$

## ۱-۶ شبیه‌سازی داکینگ مولکولی

بهبود در فرایند طراحی دارو یکی از بحث‌های چالش برانگیز در دنیای علم و تکنولوژی امروزی به‌شمار می‌رود. اصولاً در طراحی دارو از سه محیط برای بررسی ترکیب پیشنهادی استفاده می‌شود:

➤ *In silico*: بررسی ترکیب شیمیایی بالقوه دارویی با استفاده از ابزارهای نرم‌افزاری

➤ *In vitro*: بررسی ترکیب شیمیایی بالقوه دارویی در محیط کشت سلولی

➤ *In vivo*: بررسی ترکیب شیمیایی بالقوه دارویی در بدن جانوران مانند موش آزمایشگاهی

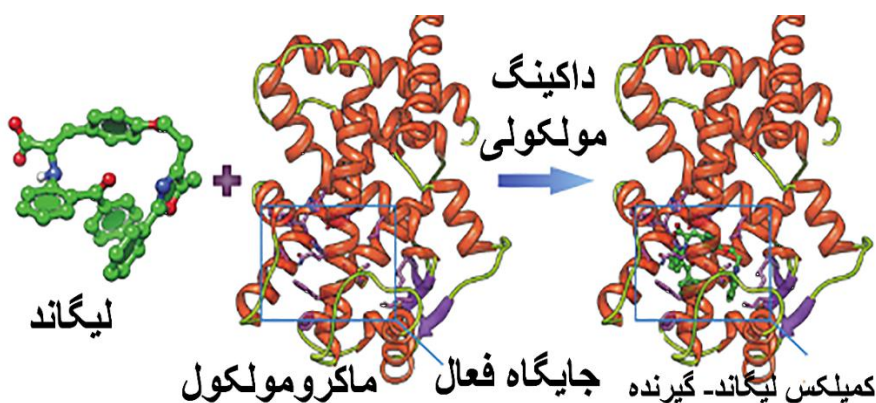
با توجه به محدودیت‌های اشاره شده در بخش مقدمه، شبیه‌سازی‌های کامپیوتری به دلیل هزینه پایین، زمان و نیروی انسانی مورد نیاز کم‌تر، جایگاه ویژه‌ای را نسبت به آنالیزهای آزمایشگاهی به خود اختصاص داده است. شبیه‌سازی دقیق و اصولی عملکرد یک ترکیب شبه دارویی با استفاده از شبیه‌سازی داکینگ مولکولی<sup>۱</sup>، دینامیک مولکولی<sup>۲</sup> و ... میسر است [۴۷، ۴۸]. از این‌رو برهم‌کنش‌های متفاوت ترکیبات شیمیایی بالقوه و ماکرومولکول (آنزیم‌ها، پروتئین‌ها و DNA) با استفاده از نرم‌افزار داکینگ مولکولی امکان‌پذیر است.

<sup>۱</sup>Molecular docking simulation

<sup>۲</sup>Molecular Dynamics (MD)



داکینگ یک الگوریتم خودکار کامپیوتری است که نحوه اتصال ترکیب به جایگاه فعال پروتئین را مشخص می‌کند. این روش شامل تعیین جهت‌گیری<sup>۱</sup> و موقعیت<sup>۲</sup> ترکیب، ساختار هندسی کنفورماسیونی و امتیازدهی می‌باشد. امتیازدهی می‌تواند معیار اندازه‌گیری انرژی اتصال، انرژی آزاد یا یک معیار عددی باشد. هر الگوریتم خودکار داکینگ به نحوی تلاش می‌کند تا ترکیب را در جهت‌گیری‌ها و کنفورماسیون‌های متفاوت در جایگاه فعال قرار دهد و امتیازی را برای هر کدام محاسبه کند. با استفاده از دانش به‌دست آمده از مطالعات داکینگ، به سنتز و بررسی ترکیبات کم‌تری نیاز خواهد شد. دلیل اصلی برای استفاده از داکینگ، پیش‌بینی ترکیباتی است که به‌خوبی به پروتئین متصل می‌شوند و علاوه بر این مشاهده ساختار هندسی ترکیب متصل شده به جایگاه فعال پروتئین است.



شکل (۱-۱) اتصال لیگاند-ماکرومولکول (گیرنده) با استفاده از مطالعه داکینگ مولکولی [۴۹]

## ۱-۷ مراحل اجرای داکینگ مولکولی

جهت اجرای داکینگ مولکولی و بررسی برهم‌کنش‌های ترکیبات فعال رعایت دو مرحله اساسی

الزامی است:

<sup>۱</sup>Orientation

<sup>۲</sup>Position

- اجرای فرایند اعتبار سنجی<sup>۱</sup> یا داک-ریداک قبل از فرایند داکینگ مولکولی که با استفاده از داکینگ لیگاند کریستالوگرافی در ساختار گیرنده در تعداد اجراهای الگوریتم ژنتیک متفاوت انجام می‌شود و داک-ریداک با کم‌ترین ریشه میانگین مربعات انحراف<sup>۲</sup> (RMSD) کم‌تر از ۱ آنگستروم)، حداقل تعداد خوشه و بیش‌ترین تعداد پیکربندی ترکیب در خوشه اول به‌عنوان شرط بهینه داکینگ انتخاب می‌شود.
- اجرای فرایند داکینگ مولکولی لیگاندهای فعال مجموعه داده‌ها و ترکیبات جدید پیشنهادی در جایگاه فعال ساختار گیرنده، در شرایط بهینه به‌دست آمده از فرایند داک-ریداک برای اجرای فرایند داکینگ و فرایند اعتبار سنجی انجام مراحل زیر ضروری است.

## ۱-۷-۱ آماده‌سازی پروتئین

صحت نتایج داکینگ به‌طور مستقیم به کیفیت ساختار کریستالوگرافی جایگاه فعال پروتئین وابسته است. ساختار کریستالوگرافی پروتئین‌ها با اشعه ایکس مشخص می‌شوند. اگر ساختار همراه با لیگاند کمپلکس شده با پروتئین و یا بدون لیگاند و کوفاکتور و فقط دارای بخش پروتئینی در جایگاه فعال فراهم شوند، به‌ترتیب ساختارهای "هولو"<sup>۳</sup> و ساختار "آپو"<sup>۴</sup> نام دارند. ساختارهای هولو مهم‌ترین روش شناسایی جایگاه فعال پروتئین هستند. مرکز ثقل لیگاندی که در پروتئین است همان مرکز تقریبی جایگاه فعال پروتئین می‌باشد. از این‌رو، ابتدا لیگاند و مولکول‌های آب موجود در ساختار کریستالوگرافی را حذف کرده

<sup>۱</sup>Validation

<sup>۲</sup>Root-mean-square deviation

<sup>۳</sup>Soaked structure

<sup>۴</sup>Apo structure

و با فرمت بانک اطلاعاتی پروتئین<sup>۱</sup> ذخیره می‌کنند و پس از ورود به محیط برنامه، هیدروژن‌هایی را که در ساختار کریستالوگرافی پروتئین دیده نمی‌شوند، اضافه می‌نمایند.

## ۱-۷-۲ ساختن لیگاند

لیگاند در برنامه داکینگ فراخوانی شده و در جایگاه فعال پروتئین هدف قرار می‌گیرد و سپس هر شبیه‌سازی داکینگ به‌طور جداگانه برای همه لیگاندهای مورد نظر توسط کاربر انجام می‌شود.

## ۱-۷-۳ تنظیم کردن جعبه شبکه‌ای<sup>۲</sup>

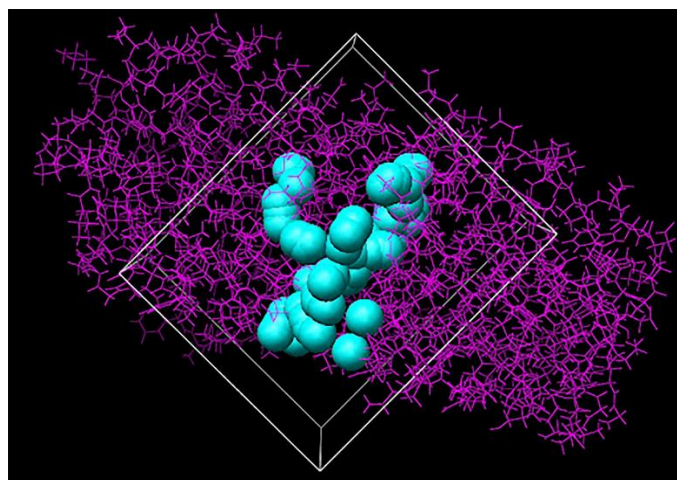
در داکینگ مولکولی نمی‌توان کل فضا را به‌عنوان جایگاه فعال گیرنده در نظر گرفت. باید بخشی از فضا که لیگاند درون جایگاه فعال قرار می‌گیرد را برای فرایند داکینگ تعریف کرد تا هم دقت بالا باشد و هم‌زمان انجام فرایند داکینگ به‌صرفه باشد. به‌منظور سرعت بخشیدن به محاسبات، یک فاصله محدود کننده<sup>۳</sup> تنظیم می‌شود، این فاصله محدود کننده معمولاً یک جعبه مستطیل شکل به نام جعبه شبکه‌ای است. بخش‌هایی از پروتئین که دور از جایگاه فعال هستند، به‌طور معمول هیچ اثر قابل اندازه‌گیری بر نتایج امتیازدهی ندارند و از این‌رو ایجاد شبکه با ابعاد یکسان برای پوشش دهی اسیدهای آمینه جایگاه فعال پروتئین توصیه می‌شود. اگر فایل مربوط به ساختار کریستالوگرافی هولو باشد به این معنی که دارای لیگاند کمپلکس شده باشد، جایگاه فعال پروتئین مشخص خواهد بود و مختصات جعبه شبکه‌ای همان مختصات جایگاه فعال و به شکل مکعبی با ابعاد یکسان تعریف می‌شود. ابعاد جعبه باید طوری انتخاب شود که اسیدهای آمینه جایگاه فعال را در بر بگیرد. شکل (۱-۲) جعبه شبکه‌ای در برگرفته اسیدهای آمینه را نشان می‌دهد [۵۰، ۵۱].

---

<sup>۱</sup>Protein data bank

<sup>۲</sup>Grid box

<sup>۳</sup>Cut off



شکل (۲-۱) جعبه شبکه‌ای [۵۲]

### ۱-۷-۴ گزینه‌های داکینگ

در هنگام تنظیم ورودی‌ها برای یک محاسبه داکینگ، گزینه‌هایی برای جایگاه فعال انعطاف‌پذیر، روش‌های نمره دهی، روش‌های جستجو، حلال پوشی، برخورد با جایگاه فعال احاطه شده و غیره موجود هستند.

### ۱-۷-۵ انجام محاسبه داکینگ

زمانی که ورودی‌ها تنظیم شدند، می‌توان محاسبات داکینگ را انجام داد. این محاسبات گاهی توسط همان کامپیوتری که صفحه رابط گرافیکی از آن استفاده شده، انجام می‌شوند و گاهی می‌توانند به یک سرور دیگر فرستاده شوند.

### ۱-۷-۶ آنالیز و تحلیل نتایج

مهم‌ترین نتیجه محاسبات داکینگ، انرژی اتصال لیگاند به جایگاه فعال است. این مقداری است که برای تعیین بهترین ترکیب مهار کننده، بین ترکیبات مختلف مقایسه می‌شود. چند حالت و وضعیت از لیگاند در جایگاه فعال که با بهترین انرژی اتصال همراه است، بررسی می‌شود تا از معقول و مناسب بودن

آن اطمینان حاصل گردد. گاهی اوقات، وضعیتی که توسط یک محاسبه داکینگ تولید می‌شود، به محقق ایده‌ای برای چگونگی تغییر ترکیبات در دوره بعدی محاسبات را می‌دهد.

دو مؤلفه کلیدی الگوریتم جستجو<sup>۱</sup> و الگوریتم رتبه‌بندی یا امتیازدهی<sup>۲</sup> در برنامه داکینگ وجود دارد. الگوریتم جستجو، مولکول را در موقعیت‌ها و صورت‌بندی‌های متفاوت در جایگاه فعال پروتئین قرار می‌دهد. انتخاب الگوریتم جستجو، مشخص می‌کند که برنامه با چه صحتی در طول فرایند، موقعیت‌های ممکن مولکول را ارزیابی می‌کند و چه مدت زمانی برای آن نیاز است. قابل ذکر است که الگوریتم جستجو صحت نتایج به دست آمده از برنامه داکینگ را مشخص نمی‌کند. الگوریتم امتیازدهی مسئول تعیین این است که آیا جهت‌گیری‌های انتخاب شده توسط الگوریتم جستجو، از لحاظ انرژی مناسب‌ترین هستند و مسئول محاسبه انرژی اتصال است.

## ۱-۷-۷ کاربردهای داکینگ مولکولی

داکینگ مولکولی کاربردهای زیادی در مطالعات مربوط به طراحی داروها دارد که شامل موارد زیر است.

۱. غربال‌گری مجازی<sup>۳</sup>

۲. بهبود طراحی و کشف دارو

۳. پیش‌بینی انرژی آزاد اتصال

۴. بررسی برهم‌کنش لیگاند-گیرنده

متداول‌ترین کاربرد داکینگ مولکولی در اتصال پروتئین-لیگاند است. هدف نهایی اتصال پروتئین-

لیگاند پیش‌بینی فعالیت بیولوژیکی لیگاند و برهم‌کنش لیگاند با پروتئین می‌باشد و در این پروژه از این

---

<sup>۱</sup>The search algorithm

<sup>۲</sup>The scoring algorithm

<sup>۳</sup>Virtual screening

هدف (اتصال لیگاند-پروتئین) جهت بررسی برهم کنش‌های ترکیبات پیشنهادی بالقوه در جایگاه فعال گیرنده استفاده می‌شود.

## ۱-۸ اهمیت مدل سازی QS(A/P)R ترکیبات شیمیایی با استفاده

از مدل‌های شبکه عصبی توسعه یافته با توصیف کننده‌های منتخب

### روش‌های جریمه‌ای و مروری بر کارهای انجام شده

با توجه به توضیحاتی که تا این بخش در مورد مراحل مدل سازی QSPR/QSAR داده شد، چالش‌هایی که یک شیمیدان محاسباتی با آن مواجه است به خوبی مشخص شد. به طور معمول، در مراحل اولیه مطالعات QSPR/QSAR، تعداد زیادی توصیف کننده مولکولی تولید می‌شود. به منظور ساخت یک مدل QSPR/QSAR با قابلیت پیش‌بینی و تفسیرپذیری بالا، یک روش انتخاب متغیر برای انتخاب زیر مجموعه بهینه از توصیف کننده‌ها الزامی است، به شرطی که مدل توسعه یافته بالاترین همبستگی بین مقادیر واقعی و پیش‌بینی شده ویژگی هدف را برای ترکیبات مورد استفاده در مدل فراهم سازد. انتخاب مؤثرترین توصیف کننده‌ها از انبوه زیادی از توصیف کننده‌های مولکولی محاسبه شده، همواره یک موضوع چالش برانگیز در ساخت مدل‌های QSPR/QSAR بوده است. در به کارگیری روش‌های انتخاب متغیر دو جنبه اصلی از جمله بهبود عملکرد پیش‌بینی برآوردگرها و ارائه تفسیر مناسبی از ارتباط بین متغیرهای مستقل و وابسته دنبال می‌شود. از این رو هدف از استفاده از روش‌های انتخاب متغیر در ساخت مدل‌های QSPR/QSAR، به دست آوردن مؤثرترین توصیف کننده‌ها با بیشترین ارتباط با متغیر وابسته است. سؤال اساسی در ارائه مدل‌های QSPR/QSAR این است که چگونه می‌توان توصیف کننده‌های مؤثر را تعیین کرد تا به طور رضایت بخشی وابستگی بین ویژگی هدف و توصیف کننده‌ها را نشان دهد. یک مدل

QSPR/QSAR با کارایی بالا، باید به ازای استفاده از حداقل پارامترهای فیزیکوشیمیایی قدرت پیش‌بینی مناسبی را در مقابل ترکیبات شیمیایی جدید نشان دهد. از این‌رو استفاده از روش‌های انتخاب متغیر کارآمد برای کاهش تعداد توصیف‌کننده‌ها توصیه می‌شود. روش‌های انتخاب متغیر با حذف توصیف‌کننده‌های زائد باعث کاهش توصیف‌کننده‌ها و انتخاب توصیف‌کننده‌های مؤثر می‌شود و در نتیجه سبب افزایش تفسیرپذیری مدل می‌شود. روش‌های کاهش ابعاد داده‌های متفاوتی در ساخت مدل‌های QSPR/QSAR کاربرد دارد. از بین روش‌های توسعه یافته، روش‌های کلاسیک به دلیل داشتن معایبی از جمله بایاس بالا در تخمین ضرایب رگرسیونی، عملکرد پایین در حضور هم‌خطی، ناپایداری در برابر تغییر نمونه‌های آزمایشی و غیره کارآمد نیستند. بنابراین استفاده از روش‌های جریمه‌شده، به دلیل رفع محدودیت‌های یاد شده در بحث مدل‌سازی QSPR/QSAR بسیار مورد توجه می‌باشند. روش‌های انتخاب متغیر جریمه شده، با انقباض ضرایب متغیرهای کم اهمیت به صفر منجر به تولید مدل‌های QSPR/QSAR تنک با قابلیت تفسیرپذیری بالا می‌شود. بنابراین استفاده از روش‌های انتخاب متغیر انقباضی به دلیل داشتن برخی مزایای ذاتی، از جمله داشتن برآوردهایی با بایاس و خطای پیش‌بینی کم و پایداری و تنگی قابل توجه در توسعه مدل‌های QSPR/QSAR شدیداً توصیه می‌شود. از این‌رو در این رساله، کارایی روش‌های انتخاب متغیر انقباضی متفاوت در مجموعه داده‌های متفاوتی مورد ارزیابی قرار گرفت.

در ادامه بحث، اهمیت ساخت مدل‌های QSPR/QSAR و تحقیقات انجام شده در سال‌های اخیر پیرامون به‌کارگیری روش‌های انتخاب متغیر انقباضی در ساخت مدل‌های QSPR/QSAR، تحت بررسی قرار گرفت.

## ۱-۸-۱ اهمیت پیش‌بینی فعالیت دارویی بازدارنده‌های ایدز با استفاده از مدل

### SCAD-ANN

ویروس نقص ایمنی اکتسابی انسانی<sup>۱</sup> (HIV) باعث ایجاد بیماری ایدز (AIDS)<sup>۲</sup> می‌شود و عامل مرگ و میر در سراسر جهان است. بر اساس داده‌های اپیدمی ایدز سازمان جهانی بهداشت<sup>۳</sup> در سال ۲۰۲۰، حدود ۱ میلیون مرگ و میر به‌اضافه ۳۷/۷ میلیون نفر مبتلا به ایدز شناخته شده است [۵۳-۵۵]. طبق نتایج منتشر شده، اگرچه مرگ و میر ناشی از ایدز در سال ۲۰۰۵ (۱/۹ میلیون) به‌طور قابل توجهی در مقایسه با سال ۲۰۱۶ (۱ میلیون) کاهش یافته است، اما متأسفانه تعداد عفونت‌های جدید در هر سال ثابت می‌ماند. یکی از دلایل اصلی کاهش مرگ و میر ناشی از ایدز، معرفی درمان ترکیبی ضد رتروویروسی (cART)<sup>۴</sup> است. با این حال، اکثر داروهای موجود در حال حاضر برای استفاده بالینی مستعد ایجاد مقاومت دارویی هستند. به‌طور خلاصه، cART چندین نقطه از چرخه همانندسازی را هدف قرار می‌دهد و معمولاً شامل بازدارنده‌های نوکلئوزیدی آنزیم نسخه‌بردار معکوس (NRTIs)، بازدارنده‌های غیر نوکلئوزیدی آنزیم نسخه‌بردار معکوس (NNRTIs) و بازدارنده‌های پروتئاز (PIs) می‌شود [۵۶].

تکثیر HIV یک فرآیند چند مرحله‌ای است و هر مرحله برای تکثیر موفقیت آمیز ویروس بسیار مهم است. به‌طور خلاصه، چرخه زندگی ویروس زمانی آغاز می‌شود که HIV به گیرنده CD4<sup>+</sup> موجود در سطح لنفوسیت T، متصل می‌شود. علاوه بر اتصال به گیرنده CD4<sup>+</sup>، HIV باید به یکی از گیرنده‌های کموکاین مانند CCR5 یا CXCR4 متصل شود تا وارد سلول شود (شکل ۱-۲). این گیرنده‌ها با پروتئین‌های پوششی ویروسی مانند گلیکوپروتئین‌های gp120 و پروتئین گذرنده gp41 تعامل دارند. هنگامی که HIV به سطح سلول میزبان نزدیک می‌شود، gp120 به گیرنده CD4<sup>+</sup> متصل می‌شود، و این امر باعث تقویت

<sup>1</sup>Human immunodeficiency virus

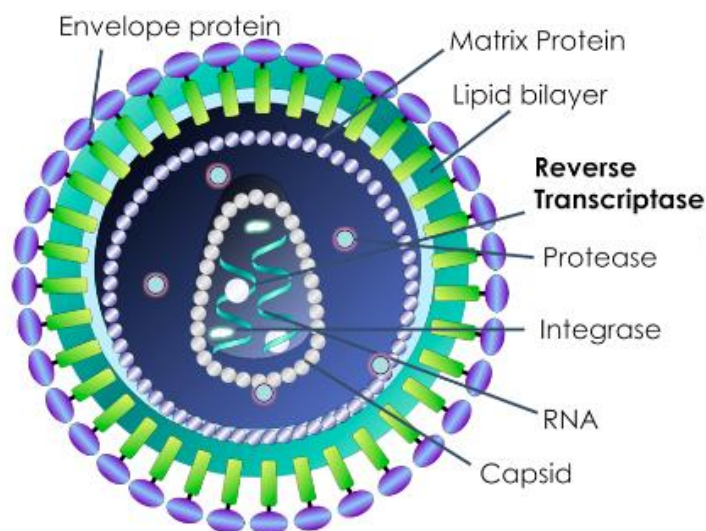
<sup>2</sup>Acquired immunodeficiency syndrome

<sup>3</sup>World Health Organization AIDS Epidemic

<sup>4</sup>Combination antiretroviral therapy



بیش تر اتصال با گیرنده مرکزی می شود و منجر به تغییر ساختاری در gp120 می شود و به gp120 اجازه می دهد تا در غشای سلول میزبان باز شود و روی خود تا شود. نوکلئوکپسید ویروس وارد سلول میزبان می شود و رشته های RNA و آنزیم های کلیدی مانند نسخه بردار معکوس (RT)، اینتگراز (IN) و پروتئاز (PR) را آزاد می کند. برای تکثیر در داخل سلول میزبان، DNA پلیمرز ویروسی از پرایمر RNA میزبان، به ویژه tRNA، برای سنتز DNA ویروسی تک رشته ای استفاده می کند، سپس با یک الگوی مشخص ویروس هیبرید RNA:DNA تشکیل می شود. نسخه بردار معکوس (RT)، رشته RNA را از هیبرید حذف می کند. سپس DNA باقی مانده ویروس به DNA ویروسی دو رشته ای تبدیل می شود. متعاقباً، اینتگراز، DNA ویروس دو رشته ای را به هسته می برد و DNA ویروس را در ژنوم سلول میزبان ادغام می کند. فعال شدن سلول باعث رونویسی DNA به mRNA می شود. سپس mRNA ویروسی به سیتوپلاسم منتقل می شود، جایی که به عنوان طرحی برای ساخت بلوک های طولانی تر از ویروس جدید استفاده می شود [۵۴، ۵۷].



شکل ۲-۱ ساختار کلی HIV [۵۴]

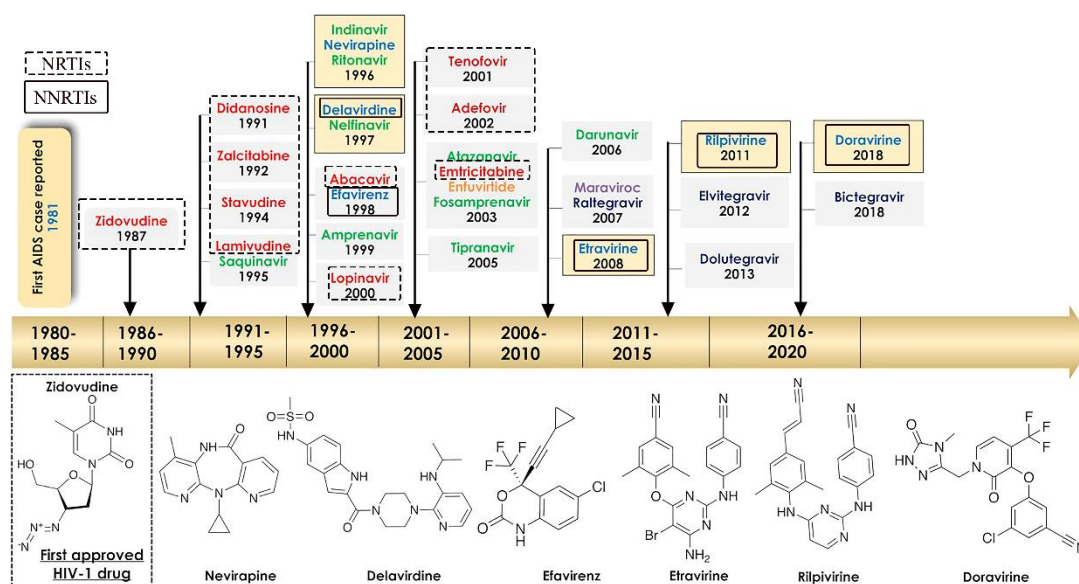
داروهایی که در مراحل کلیدی تکثیر تداخل ایجاد می کنند، باعث توقف آن ها در مراحل زیر

می شوند [۵۸].

- مسدود کردن ورود ویروس به سلول میزبان با مهارکننده‌های ورودی، از جمله مهارکننده‌های همجوشی<sup>۱</sup> (FIs)
  - متوقف نمودن تکثیر با مهارکننده‌های نوکلئوزیدی نسخه بردار معکوس (NRTIs) و یا مهارکننده‌های نسخه بردار معکوس غیر نوکلئوزیدی (NNRTIs)
  - مهار ادغام DNA ویروسی با ماده ژنتیکی میزبان با مسدود کردن فعالیت اینتگرز (InSTI).
  - انسداد بلوغ بیش‌تر ویروس با مهارکننده پروتئاز (PI)
- رونویسی معکوس یک DNA پلیمرز وابسته به RNA است که از یک رشته RNA برای سنتز DNA ویروسی دو رشته‌ای استفاده می‌کند [۵۹]. بر اساس محل اتصال و طبقه شیمیایی ترکیبات، مهارکننده‌های RT به NRTIs و NNRTIs طبقه‌بندی می‌شوند [۶۰، ۶۱]. شکل ۱-۳ نقطه عطف کوتاهی در توسعه داروهای HIV را نشان می‌دهد و همان‌طور که مشخص است، مهارکننده‌های مهم RT دارای برجستگی قابل توجهی می‌باشند. NNRTIs (nevirapine، delavirdine، efavirenz، etravirine، rilpivirine و doravirine) از دسته داروهای RT مورد تأیید سازمان غذا و داروی ایالات متحده<sup>۲</sup> است. ترکیبات NNRTIs با قدرت بیش‌تر و سمیت کم‌تر همواره برای مقابله با HIV بسیار کارآمد هستند. از طرفی، با توجه به مقاومت دارویی سریع سویه‌های ویروسی جدید، پیشنهاد NNRTI ها با اثربخشی بهتر بسیار حائز اهمیت است [۵۴].

<sup>1</sup>Fusion Inhibitor

<sup>2</sup>U.S. Food and Drug Administration



شکل ۳-۱ جدول زمانی کوتاهی از توسعه داروهای HIV، داروهای شناسایی شده برای RT با خط چین (NRTIs) و خط جدا شده اند [۵۴] (NNRTIs)

بنابراین با توجه به اهمیت بازدارنده‌های غیرنوکلئوزیدی نسخه بردار معکوس در مقابله با HIV، پیشنهاد ترکیبات جدید NNRTIs با قدرت بیش‌تر، سمیت کم‌تر در امر طراحی دارو حائز اهمیت است. بنابراین در اولین بخش مطالعات تجربی این رساله تلاش شد و یک مدل QSAR برای پیش‌بینی فعالیت دارویی بازدارنده‌های NNRTIs توسعه یافت. همان‌طور که در بخش ۱-۸ توضیح داده شد، انتخاب مؤثرترین توصیف‌کننده‌ها موضوع چالشی در ساخت مدل‌های QSAR است و استفاده از روش‌های انتخاب متغیر انقباضی به دلیل داشتن مزایای ذاتی متفاوت، حائز اهمیت است. بنابراین کارایی روش انتخاب متغیر انقباضی SCAD در ساخت مدل QSAR مورد بررسی قرار گرفت. بنابراین در ادامه و در بخش مروری بر کارها، به بررسی مطالعات و پروژه‌هایی که در سال‌های اخیر در مورد به‌کارگیری روش انتخاب متغیر SCAD در ساخت مدل‌های QSAR بحث شده است پرداخته خواهد شد.

## ۱-۸-۲ اهمیت پیش‌بینی فعالیت دارویی بازدارنده‌های SARS-CoV-2 با استفاده

### از مدل ALASSO-ANN

سندرم حاد تنفسی ویروس کرونا<sup>۱</sup>(SARS-CoV-2)، به‌عنوان هفتمین ویروس کرونا‌ی انسانی، در ووهان، استان هوبی، چین، در طی اپیدمی اخیر ذات‌الریه در ژانویه ۲۰۲۰ کشف شد [۶۲، ۶۳]. از آن زمان، این ویروس در سراسر جهان گسترش یافته است و تا به الان، حدود ۲۵۸ میلیون عفونت و حدود ۵ میلیون قربانی داشته است [۶۴]. طبق آمار منتشر شده، SARS-CoV، SARS-CoV-2 و کروناویروس سندرم تنفسی خاورمیانه (MERS-CoV)<sup>۲</sup> به‌ترتیب باعث ذات‌الریه شدید با میزان مرگ و میر ۲/۹٪، ۹/۶٪ و ~۳۶٪ می‌شوند [۶۵-۶۷].

SARS-CoV-2 متعلق به طبقه بتا-کرونا ویروس از خانواده کروناویردایی<sup>۳</sup> و از راسته نیدوویرالز<sup>۴</sup> است. SARS-CoV-2 بسیار شبیه به کروناویروس‌های SARS مانند است [۶۲] و تقریباً ۹۶٪ دارای اشتراک ساختاری هستند. پروتئین‌های SARS-CoV-2 دو سوم ژنوم ویروس را اشغال می‌کنند و اهداف برجسته‌ای برای داروهای تعدیل‌کننده تکثیر ویروس هستند. این گروه شامل آنزیم‌های تکثیر و رونویسی ویروسی RNA پلیمراز، هلیکاز و پروتئاز شبه پاپائین (PL<sup>Pro</sup>)<sup>۵</sup> همچون پروتئاز ۳C-مانند (۳CL<sup>Pro</sup> و یا M<sup>Pro</sup>) می‌شود [۶۸]. SARS-CoV-2 به‌عنوان یک ویروس RNA تک‌رشته‌ای، مستقیماً پس از عفونت، سلول از سیستم‌های ترجمه سلولی برای تولید پلی‌پروتئین‌های ویروسی استفاده می‌کند [۶۹]. غربالگری کتابخانه‌ای ترکیبات و مطالعات داکینگ مولکولی بر اساس ساختار کریستالوگرافی SARS-CoV و MERS، مهارکننده‌های 3CL<sup>Pro</sup> متفاوت، عملکرد بالقوه‌ای را علیه SARS-CoV-2 نشان داده است [۷۰، ۷۱].

<sup>1</sup>Severe acute respiratory syndrome coronavirus

<sup>2</sup>Middle East respiratory syndrome coronavirus

<sup>3</sup>Coronaviridae

<sup>4</sup>Nidovirales

<sup>5</sup>Papain-like protease

بنابراین پیشنهاد ترکیبات  ${}^3\text{CL}^{\text{pro}}$  جدید با فعالیت دارویی بالقوه از اهمیت بالایی برخوردار است. از این رو ارائه یک مدل QSAR قابل اعتماد برای پیش‌بینی فعالیت دارویی ترکیبات  ${}^3\text{CL}^{\text{pro}}$  جدید توصیه می‌شود. در دومین بخش این رساله، یک مدل QSAR با استفاده از توصیف‌کننده‌های منتخب روش انتخاب متغیر انقباضی لاسو تطبیقی (ALASSO) برای پیش‌بینی فعالیت دارویی بازدارنده‌های  ${}^3\text{CL}^{\text{pro}}$  توسعه یافت.

## ۱-۸-۳ اهمیت پیش‌بینی فعالیت دارویی برخی از بازدارنده‌های ایدز و سرطان با

### استفاده از مدل LAD-LASSO-ANN

عادات غذایی نامناسب و سبک زندگی وابسته به میز، باعث بروز بیماری‌های مختلفی مانند سرطان، عوارض قلبی عروقی، دیابت و اختلالات ایمنی می‌شود [۷۲]. سرطان یک مشکل عمده بهداشت عمومی در سراسر جهان است و دومین عامل مرگ و میر در دنیا است. متأسفانه نرخ رشد آن در کشورهای توسعه یافته‌تر و در حال توسعه دیده می‌شود. در سال ۲۰۲۰، تشخیص و درمان سرطان به دلیل همه‌گیری بیماری کروناویروس ۲۰۱۹ با مشکل مواجه شد. به‌عنوان مثال، کاهش دسترسی به مراقبت به دلیل بسته شدن مراکز مراقبت‌های بهداشتی منجر به تأخیر در تشخیص و درمان می‌شود، که ممکن است منجر به کاهش کوتاه مدت در بروز سرطان و به دنبال آن افزایش در مرحله پیشرفته بیماری و در نهایت افزایش مرگ و میر شود. از این رو، ارائه دقیقی از آمار مبتلایان و مرگ و میر دقیق زمان‌بر است [۷۳]. طبق اعلامیه سازمان جهانی بهداشت در سال ۲۰۲۰، سرطان حدود ۱۰ میلیون مرگ و میر را به دنبال داشته است و طبق تقسیم‌بندی انجام شده، سرطان سینه با ۲/۲۶ میلیون جمعیت، در رتبه اول مرگ و میر می‌باشد. سرطان ریه و سرطان روده به ترتیب با ۲/۲۱ و ۱/۹۳ میلیون مرگ و میر رتبه دوم و سوم سرطان پرخطر را به خود اختصاص داده‌اند. در ادامه سرطان پروستات، پوست و معده مجموعاً ۳/۵ میلیون مرگ و میر را به همراه داشته‌اند [۶۴]. با توجه به آمار اشاره شده حدود ۲۲ درصد و ۱۹ درصد از کل مرگ و میر مربوط به بیماری

سرطان به ترتیب به سرطان‌های ریه و روده تعلق دارد. نرخ بقای بیماران مبتلا به سرطان ریه و روده، علی‌رغم پیشرفت‌هایی که در درمان دارویی و تشخیص آن صورت گرفته است، هنوز بسیار پایین است بنابراین، تحقیق برای دستیابی به یک رویکرد درمانی مناسب ضروری است [۷۴].

فسفوئینوزیتید-۳-کینازها<sup>۱</sup> (PI3Ks)، کینازهای لیپیدی هستند که در پاسخ به فاکتورهای رشد توسط تعداد معینی از پروتئین گیرنده تیروزین کیناز (RTKs)<sup>۲</sup>، یعنی گیرنده فاکتور رشد اپیدرمی (EGFR)<sup>۳</sup>، گیرنده فاکتور رشد اپیدرمی انسانی (HER2)<sup>۴</sup> حساس به فعال شدن هستند. PI3K ها به سه طبقه اصلی تقسیم می‌شوند [۷۵]. شواهد نشان می‌دهد جهش فعال کننده PI3K $\alpha$  معمولاً در بسیاری از سرطان‌ها مشاهده می‌شود. بنابراین، مهارکننده‌های مولکولی کوچک علیه PI3K $\alpha$  برای درمان سرطان حائز اهمیت است [۷۶]. تحقیقات زیادی در مورد مهار کننده‌های PI3K $\alpha$  وجود دارد و برخی از مهارکننده‌ها وارد آزمایش‌های بالینی نیز شده‌اند. چندین ماده شیمیایی، از جمله بازدارنده‌های نسل اول PI3K مانند ورتمانین<sup>۵</sup> (شکل ۱-۴ A)، به‌عنوان مهارکننده‌های PI3K به‌طور گسترده شناخته شده‌اند. مهارکننده‌های نسل دوم PI3K مولکول‌های کوچکی هستند که از جمله می‌توان به LY294002 (شکل ۱-۴ B) به‌عنوان آنالوگ کوئرستین<sup>۶</sup> (شکل ۲-۱ C) اشاره کرد [۷۷]. مطالعات سرطان شناسی، LY294002 و آنالوگ‌های آن را به‌عنوان قوی‌ترین مهارکننده‌های PI3K معرفی نموده است [۷۸، ۷۹]. طبق تحقیقات انجام شده مشتقات کرومنو [۳، ۴-c] پیرازول-۴(۲H)-اون<sup>۷</sup> حاوی زنجیره بنزیل و آلکیل به‌عنوان مهار کننده‌های بالقوه PI3K $\alpha$  و از دسته آنالوگ‌های B شناخته شده هستند [۸۰]. این مشتقات در برابر سلول‌های سرطانی ریه

<sup>1</sup>Phosphoinositide-3-kinases

<sup>2</sup>Receptor protein tyrosine kinases

<sup>3</sup>Epidermal growth factor receptor

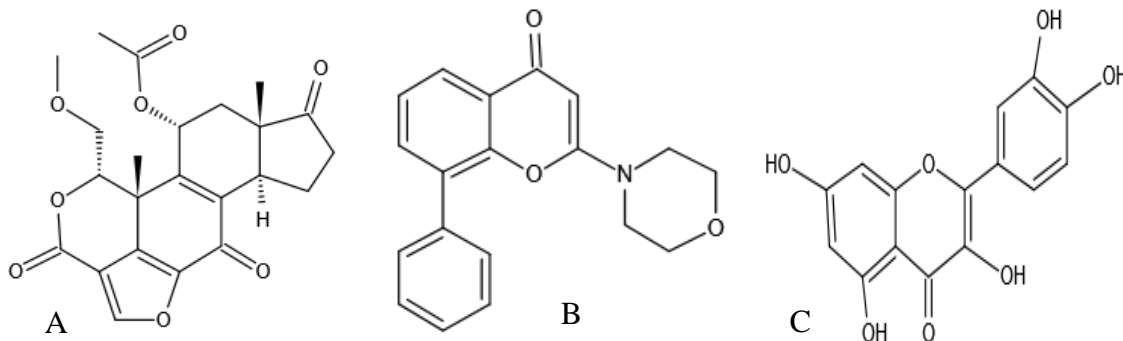
<sup>4</sup>Human epidermal growth factor receptor 2

<sup>5</sup>Wortmannin

<sup>6</sup>Quercetin

<sup>7</sup>Chromeno[4,3-c]pyrazol-4(2H)-one

(A549) و کولورکتال (HCT-116) مورد آزمایش قرار گرفته و فعالیت دارویی مناسب و اثرات ضد تکثیری قابل توجهی را نشان دادند [۸۱].



شکل ۴-۱ شکل بازدارنده‌های بالقوه PI3K. A: ورتمانین، B: LY294002 و C: کوئرستین

بنابراین طبق توضیحات ارائه شده در مورد اهمیت بازدارنده‌های PI3K، پیش‌بینی فعالیت ترکیبات آنالوگ و پیشنهاد ترکیبات جدید سنتز نشده از اهمیت بسیاری برخوردار است. بنابراین در این مطالعه از اطلاعات ساختاری و فعالیت دارویی ضد سرطانی روده و ریه این مشتقات استفاده شد تا مدل‌های QSAR با قدرت پیش‌بینی قابل اعتماد معرفی شود. به‌منظور افزایش قدرت پیش‌بینی و تفسیرپذیری مدل، استفاده از یک روش انتخاب متغیر انقباضی کارآمد توصیه می‌شود. از این‌رو در سومین بخش این رساله از روش انتخاب متغیر انقباضی حداقل انحراف مطلق - حداقل قدر مطلق انقباض و عملگر انتخاب کننده (LAD-<sup>۱</sup> LASSO) به دلیل مزایایی چون تنگی، مقاوم به وجود داده‌های دور افتاده و نزدیک بودن به یک مدل حقیقی (اوراکل) [۳۲] برای کاهش ابعاد داده‌ها قبل از مدل‌سازی با ANN استفاده شد. در ادامه مطالعات انجام شده در زمینه QSAR برای ساخت مدل‌های طبقه‌بندی و رگرسیونی مورد بررسی قرار گرفته است.

<sup>۱</sup>Least absolute deviation (LAD)- least absolute shrinkage and selection operator (LASSO)

## ۱-۸-۴ اهمیت پیش‌بینی شاخص بازداری ترکیبات آلی فرار با استفاده از مدل

### SCAD-ANN

ترکیبات آلی فرار (VOCs)<sup>۱</sup> مواد شیمیایی آلی هستند که به‌طور کلی دارای فشار بخار بالا و حلالیت پایین در آب هستند. ترکیبات مختلف زیادی وجود دارد که ممکن است به‌عنوان ترکیبات آلی فرار طبقه‌بندی شوند. بسیاری از VOC ها مواد شیمیایی ساخته شده توسط انسان هستند که در ساخت رنگ‌ها، مواد دارویی و خنک‌کننده‌ها استفاده و تولید می‌شوند. VOC ها آلاینده‌های رایج آب‌های زیرزمینی یا آلاینده‌های خطرناک هوا هستند [۸۲]. VOC ها همچنین نقش عمده‌ای در تشکیل آلاینده‌های ثانویه مختلف از طریق واکنش‌های فیتوشیمیایی در حضور نور خورشید و اکسیدهای نیتروژن دارند. علاوه بر این، برخی از VOC ها می‌توانند در تخریب لایه اوزون و ایجاد آلودگی‌های پایدار در مناطق نقش داشته باشد. بنابراین، این ترکیبات در طول دو دهه اخیر یک موضوع مهم زیست محیطی بوده و توجه گروه‌های مختلف تحقیقاتی را به خود جلب کرده است [۸۳-۸۸]. روش‌های تجزیه‌ای ویژه با حساسیت و گزینش پذیری بالا برای شناسایی و اندازه‌گیری VOC ها در اختیار دانشمندان قرار دارد که از جمله آن‌ها می‌توان به استفاده از کروماتوگرافی گازی در شناسایی و اندازه‌گیری شاخص بازداری این ترکیبات اشاره کرد [۸۸]. محدودیت‌هایی برای استفاده از این روش آزمایشگاهی وجود دارد. به‌طوری‌که در برخی موارد استانداردهای موجود با درجه خلوص بالا وجود ندارد یا در برخی موارد ممکن است این استانداردهای مورد نیاز در دسترس نباشند. عدم دسترسی به ستون‌های متفاوت جهت انجام آنالیز و یا عدم تکرارپذیری داده‌های اندازه‌گیری شده، همگی، از جمله محدودیت‌هایی است که به‌کارگیری روش‌های تئوری در پیش‌بینی شاخص بازداری این ترکیبات را حائز اهمیت ویژه‌ای می‌سازد [۸۹-۹۱].

---

<sup>۱</sup>Volatile organic compounds



یک جایگزین معتبر برای مقابله با این محدودیت‌ها، ادغام اطلاعات ساختاری VOC های شناسایی شده و شاخص‌های بازداری تجربی آن‌ها به منظور ارائه یک مدل نظری برای تخمین RI مربوط به VOC های جدید است. بنابراین ارتباط کمی ساختار - ویژگی (QSPR)، یک روش امیدوارکننده‌ای را برای تخمین شاخص بازداری بر اساس توصیف‌کننده‌های ساختار مولکولی فراهم می‌سازد [۸۳، ۸۴، ۹۴-۹۲].

## ۱-۹ مروری بر تحقیقات انجام شده در مورد به کارگیری روش‌های انتخاب

### متغیر جریمه‌ای در ساخت مدل‌های QSAR/QSPR

در این بخش جستجوی کتابخانه در مورد به کارگیری روش‌های انقباضی در ساخت مدل‌های QSAR/QSPR انجام شده است که به ترتیب سال در ادامه مرتب شده‌اند.

با توجه به جستجوی انجام شده در پایگاه داده‌های متفاوت، از روش‌های انقباضی SCAD، ALASSO و LAD-LASSO در مطالعات QSAR/QSPR به دو شکل طبقه‌بندی<sup>۱</sup> و رگرسیونی استفاده شده است. موارد مربوط به کاربرد روش‌های انقباضی در مطالعات طبقه‌بندی<sup>۲</sup> QSAR/QSPR مد نظر اهداف این رساله نبوده است و جزییات آن در بخش مرور بر کارها ذکر نشده است [۱۰۶-۹۵]. لازم به ذکر است که در برخی از این تحقیقات، روش‌های انقباضی SCAD، ALASSO و LAD-LASSO تنها به منظور مقایسه با روش پیشنهادی انقباضی مورد نظر نویسندگان، به کار گرفته شده است و لزوماً به عنوان روش برتر مطالعه برگزیده نشده است. در ادامه به مرور مطالعات منتشر شده از به کارگیری روش‌های انقباضی SCAD، ALASSO و LAD-LASSO در مدل‌های رگرسیونی QSAR/QSPR پرداخته خواهد شد.

---

<sup>۱</sup>Classification

<sup>۲</sup>Classification

پنگ<sup>۱</sup> و همکارانش در سال ۲۰۰۶ نقطه جوش ۵۳۰ هیدروکربن را با استفاده از جفت مدل SCAD و مدل Kriging پیش‌بینی نمودند. مدل SCAD-Kriging با ۱۵ توصیف‌کننده با مدل حداقل مربعات خطی معمولی (OLS)<sup>۲</sup> مقایسه شد. مقدار ریشه میانگین مربعات خطا (RMSE)<sup>۳</sup> برای مدل پیشنهادی نویسنده ۱۷٪ کاهش یافت [۱۰۷].

الجمال و همکارانش در سال ۲۰۱۵ مدل QSAR را برای پیش‌بینی فعالیت دارویی (pIC<sub>50</sub>)<sup>۴</sup> مشتق ایمیدازو [b-۴،۵] پیریدین به‌عنوان بازدارنده‌های سرطان توسعه دادند. در این مطالعه از فرم اصلاح شده‌ای از روش انقباضی ALASSO به‌عنوان مدل استفاده شد و کارایی روش با مدل‌های LASSO و ALASSO مقایسه شد. به‌طوری‌که روش‌های انتخاب متغیر انقباضی بر روی کل توصیف‌کننده‌های محاسبه شده، ۲۵ بار (با تقسیم‌بندی مجموعه آموزش و آزمون متفاوت) اجرا شد و توصیف‌کننده‌ها بر اساس فراوانی ۵۰٪ به بالای تکرار توصیف‌کننده‌ها، برای مدل‌سازی مورد استفاده قرار گرفتند. پارامتر  $Q^2$  برای داده‌های مجموعه آزمون به‌ترتیب برای مدل‌های QSAR توسعه یافته با مدل‌های انقباضی ALASSO اصلاح شده (۲۶ متغیر)، LASSO (۱۶ متغیر) و ALASSO (۱۲ متغیر) برابر با ۰/۸۷، ۰/۸۰ و ۰/۷۶ به‌دست آمد. نتایج نشان می‌دهد که مدل توسعه یافته انقباضی دارای قدرت پیش‌بینی مناسبی هستند [۱۰۸].

الجمال<sup>۴</sup> و همکارانش در سال ۲۰۱۶، روش انقباضی Bridge با نُرم  $L_{1/2}$  را برای ساخت مدل QSAR ۱۱۱ بازدارنده‌های استیل کولین استراز<sup>۵</sup> ارائه کردند. به‌منظور مقایسه روش پیشنهادی از روش‌های انتخاب متغیر دیگری همچون LASSO، ALASSO و SCAD نیز استفاده شد. پارامترهای آماری متفاوتی برای این مدل‌ها محاسبه شد. پارامتر ضرایب تعیین مجموعه آزمون به ترتیب برای مدل‌های Bridge،

---

<sup>1</sup>Peng

<sup>2</sup>Ordinary least squares

<sup>3</sup>Root mean square error

<sup>4</sup>Algamal

<sup>5</sup>Acetylcholinesterase

<sup>6</sup>Least Absolute Shrinkage and Selection Operator

<sup>7</sup>Adaptive least Absolute Shrinkage and Selection Operator

LASSO, ALASSO و SCAD به ترتیب برابر با ۰/۹۱، ۰/۸۰، ۰/۸۱ و ۰/۸۷ به دست آمد. نتایج نشان می‌دهد که همه روش‌های انتخاب متغیر به کار گرفته شده در این مطالعه دارای قدرت پیش‌بینی قابل قبولی هستند [۱۰۹].

الجمال و همکارانش در سال ۲۰۱۷، مدل‌های اصلاح شده تنک مبتنی بر ضرایب تنظیم‌کننده متفاوت (RS-QSPR) و هم‌چنین مدل SCAD را برای پیش‌بینی شاخص بازداري ۱۶۹ ترکیب اسانس ارائه کردند. مدل‌های RS-QSPR و SCAD به ترتیب ۴ و ۹ توصیف‌کننده مؤثر را انتخاب کردند. ضریب تعیین مجموعه آزمون برای این دو مدل به ترتیب برابر با ۰/۸۷ و ۰/۹۴ به دست آمد. نتایج نشان‌دهنده قدرت پیش‌بینی هر دو مدل QSPR توسعه یافته می‌باشد [۱۱۰].

الجمال و همکارانش در سال ۲۰۱۸ مدل‌های انقباضی QSAR (LASSO, Bridge و PBridge) را برای پیش‌بینی نقطه ذوب ۶۰ ترکیب نیتروآروماتیک ارائه دادند. توصیف‌کننده‌های منتخب هر روش به ترتیب برابر با ۶، ۷ و ۴ برای روش‌های BRIDGE, LASSO و PBridge به دست آمد. پارامترهای آماری متفاوتی برای مدل‌های توسعه یافته محاسبه شد. پارامتر ضریب تعیین مجموعه داده‌های خارجی برای مدل‌های انقباضی یاد شده به ترتیب برابر با ۰/۹۱، ۰/۸۱ و ۰/۹۴ به دست آمد. نتایج حاکی از اعتبار قابل قبول مدل‌های انقباضی توسعه یافته است [۱۱۰].

ماجومدار<sup>۲</sup> و همکارانش در سال ۲۰۱۸ مدل‌های QSAR را برای مجموعه داده همگن متشکل از ۹۵ آمین، ارائه نمودند. با استفاده از توصیف‌کننده‌های محاسبه‌شده، فعالیت جهش‌زایی این ترکیبات را با استفاده از روش‌های مدل‌سازی متفاوتی همچون  $PLS^3$ ،  $RF^4$ ، LASSO و SCAD پیش‌بینی کردند. میانگین مربعات خطای پیش‌بینی برای مدل‌های QSAR مربوطه به دست آمد و به ترتیب برابر با ۲۹/۱۱،

---

<sup>1</sup>robust sparse QSPR

<sup>2</sup>Majumdar

<sup>3</sup>Principal component regression

<sup>4</sup>Partial least squares regression

<sup>5</sup>Random forest

۱۸/۸۶، ۱۷/۲۵، ۲۶/۸۵، ۲۵/۸۱ است. نتایج نشان می‌دهد که مدل‌های ساخته شده از اعتبار مناسبی برخوردار هستند [۱۱۱].

یونگ زی<sup>۱</sup> و همکارانش در سال ۲۰۱۸، چندین روش انقباضی از جمله SCAD، MCP، EN<sup>۲</sup> را برای پیش‌بینی فعالیت دارویی ۴ دسته داده عمومی به کار گرفتند. پارامترهای آماری متفاوتی برای ارزیابی مدل‌های انقباضی توسعه یافته محاسبه شدند. نتایج مربوط به هر ۴ دسته داده از اعتبار قابل قبولی برخوردار بودند [۱۱۲].

علاوه بر این طبق جستجوی کتابخانه‌ای انجام شده، در سال‌های اخیر گزارشی مبنی بر استفاده از روش LAD-LASSO در پیش‌بینی فعالیت دارویی ترکیبات متفاوت در قالب ارائه یک مدل QSAR وجود نداشته است. تنها یک گزارش از به‌کارگیری روش‌های مدل‌سازی انقباضی LAD-Bridge، LAD-SCAD و LAD-LASSO در طبقه‌بندی ترکیبات فعال و غیرفعال ۱۰۸ بازدارنده آنفلوانزا توسط الجمال و همکارش در سال ۲۰۱۹ منتشر شده است [۱۱۳]. بنابراین با توجه به این که این روش در مطالعات QSAR استفاده نشده است، در این رساله برای اولین بار از روش انتخاب متغیر LAD-LASSO جفت شده با ANN برای پیش‌بینی فعالیت دارویی بازدارنده‌های سرطان ریه، روده و بازدارنده‌های ایدز استفاده شده است.

---

<sup>۱</sup>Yong Xia

<sup>۲</sup>Minimax concave penalty

<sup>۳</sup>Elastic-net

## ۱- ۱۰ نوآوری تحقیق

با توجه به موارد اشاره شده در بخش مروری بر کارها، به خوبی مشخص شد که به کارگیری روش‌های انقباضی به دلیل داشتن مزایای ذاتی چون تنکی، پایداری، بایاس کم و کارایی بالا در انتخاب توصیف کننده‌های مؤثر در سال‌های اخیر مورد توجه محققین بوده است. به طوری که مطالعات متفاوتی مبنی بر به کارگیری روش‌های انقباضی برای ساخت مدل QSAR/QSPR ارائه شده است. نتایج جستجوی کتابخانه‌ای نشان می‌دهد که مطالعات اندکی به بررسی کارایی روش‌های انقباضی به عنوان روش انتخاب متغیر مناسب و ترکیب آن با روش‌های مدل‌سازی خطی و غیر خطی پرداخته‌اند. به این منظور، در این رساله، با توجه به کارایی و قدرت پیش بینی مدل غیر خطی شبکه عصبی مصنوعی (ANN)؛<sup>۱</sup> مدل‌های هیبریدی مبتنی بر روش‌های انقباضی به عنوان روش انتخاب متغیر تنک و پایدار و مدل شبکه عصبی به عنوان مدل غیر خطی قدرتمند در ساخت مدل‌های QSAR/QSPR پیشگو برای پیش بینی فعالیت دارویی/ویژگی ترکیبات شیمیایی ارائه شد [۱۱۴، ۱۱۵].

بنابراین هدف اصلی این رساله ایجاد مدل‌های QSAR تنک و تفسیرپذیر با حداقل تعداد توصیف کننده مناسب برای پیش بینی فعالیت دارویی/ویژگی ترکیبات شیمیایی است. به این منظور با توجه به کارایی و پراکندگی روش‌های انقباضی هم‌چون SCAD، ALASSO و LAD-LASSO، در انتخاب توصیف کننده‌هایی با بیشترین تأثیر بر متغیر وابسته، تلاش شد تا به صورت هدفمندی از این روش‌ها در ساخت مدل‌های QSAR/QSPR مبتنی بر شبکه عصبی مصنوعی استفاده شود. لازم به ذکر است پس از ارزیابی موفق مدل‌های QSAR/QSPR توسعه یافته، از مدل‌های شبکه عصبی توسعه یافته با توصیف کننده‌های منتخب روش‌های انقباضی متفاوت، برای پیشنهاد ترکیبات جدید بالقوه استفاده شده است.

---

<sup>۱</sup>Artificial Neural Network

ش تخریبی



## ۲-۱ معرفی نرم افزارهای مورد استفاده برای مدل سازی QSAR

در این رساله از نرم افزارهای متفاوتی برای مدل سازی استفاده شده است که در ادامه به طور مختصر به توضیح و بررسی هر یک از آنها پرداخته شده است.

### ۲-۱-۱ نرم افزار هایپرکم

نرم افزار هایپرکم [۱۱۶] به عنوان یک ابزار پیشرفته در شیمی محاسباتی، برای رسم و بهینه سازی مولکول های ساده و پیچیده با استفاده از روش های مکانیکی و کوانتومی شناخته می شود. با استفاده از این نرم افزار می توان، طول پیوند، زاویه پیوند، زاویه پیچشی و غیره ترکیبات شیمیایی را بهینه کرد. هر چند برخی از توصیف کننده های ترکیبات شیمیایی از قبیل ضریب شکست، قطبش پذیری، چربی دوستی و غیره با نرم افزار هایپرکم قابل محاسبه است، ولی غالباً برای محاسبه توصیف کننده های ترکیبات استفاده نمی شود و در مطالعات QSAR، خروجی این نرم افزار به عنوان ورودی برای نرم افزار دراگون به کار گرفته می شود.

### ۲-۱-۲ نرم افزار دراگون<sup>۱</sup>

از نرم افزار دراگون جهت استخراج توصیف کننده های مولکولی استفاده می شود. نرم افزار دراگون توسط گروه کمومتریکس دانشگاه میلانو طراحی شده است [۱۱۷]. نرم افزار دراگون نسخه ۵/۵، قادر به محاسبه ۳۲۲۴ توصیف کننده می باشد که به ۲۲ دسته اصلی تقسیم می شوند. توصیف کننده هایی از جمله نوع و تعداد اتم ها، گروه های عاملی، توصیف کننده های توپولوژیکی و هندسی و غیره با استفاده از این نرم افزار قابل محاسبه است. اطلاعات کاملی از جمله ماهیت توصیف کننده ها، فرمول محاسبه، مراجع و جزییات توصیف کننده ها در کتاب مرجع توصیف کننده های مولکولی قابل دسترس است [۹].

---

<sup>1</sup>Dragon



## ۲-۱-۳ نرم افزار SPSS<sup>۱</sup>

SPSS یکی از توانمندترین و جامع ترین نرم افزارهای آماری برای تحلیل داده ها است. نرم افزار آماری SPSS ساخته دانشگاه استنفورد<sup>۲</sup> در سال ۱۹۷۰ است. از جمله کاربردهای نرم افزار آماری SPSS می توان به موارد زیر اشاره کرد:

- ❖ تحلیل و آنالیز داده های ورودی
- ❖ تهیه جداول و نمودارهای آماری
- ❖ به دست آوردن توزیع و رفتار داده های ورودی
- ❖ ایجاد داده های تصادفی
- ❖ پردازش انواع رگرسیون
- ❖ محاسبه انواع آزمون های آماری

در این رساله برای انجام آنالیزهای آماری متفاوت روی توصیف کننده های استخراج شده نرم افزار دراگون، اجرای روش رگرسیون خطی چندگانه برای انتخاب مؤثرترین توصیف کننده ها، محاسبه هم خطی<sup>۳</sup> بین توصیف کننده ها از نرم افزار SPSS 25 استفاده شد [۱۱۸].

## ۲-۱-۴ نرم افزارهای R و R-studio

نرم افزار R، یک زبان برنامه نویسی و محیط نرم افزاری رایگان برای محاسبات آماری و علم داده ها است. نرم افزار R دارای بسته های آماری و تحلیلی مختلف بر اساس موضوع مورد نیاز است، در نتیجه بسیار مورد توجه محققین قرار گرفته است. از جمله کاربردهای نرم افزار R می توان به موارد زیر اشاره کرد:

➤ برنامه نویسی و محیط نرم افزاری برای محاسبات آماری و علم داده ها.

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

<sup>۲</sup>Stanford University

<sup>۳</sup>Collinearity

<sup>۴</sup>Packages

- برنامه‌نویسی ساده و پیشرفته شامل عبارتهای شرطی، حلقه و غیره.
  - انجام عملیات داده‌کاوی و یادگیری ماشین مانند دسته‌بندی، خوشه‌بندی، انتخاب متغیر، مدل‌سازی، تحلیل شبکه و غیره با استفاده از کتابخانه‌های نرم افزار.
  - امکان ذخیره، بازیابی و دست‌کاری داده‌ها.
- R-Studio نیز یک ابزار رایگان برای استفاده از نرم‌افزار R است و امکان توسعه و به اشتراک‌گذاری آنالیزهای آماری را برای محققین فراهم می‌کند. نرم‌افزار R-Studio دارای محیط گرافیکی و نمایشی بهتری نسبت به نرم‌افزار R است و R-Studio، نرم‌افزار کاربرپسند<sup>۱</sup> بوده و استفاده از آن، برای اجرای بسته‌های نرم‌افزاری آماری موجود در نرم‌افزار R ضروری می‌باشد.
- در این رساله از نرم‌افزارهای R و R-studio برای پیش‌پردازش داده‌ها، تقسیم‌بندی داده‌ها، اجرای روش‌های رگرسیون انقباضی برای انتخاب مؤثرترین توصیف‌کننده‌ها و هم‌چنین، اجرای آزمون‌های متفاوت ارزیابی مدل‌های برتر استفاده شد.

## ۲-۱-۵ نرم‌افزار متلب

متلب<sup>۲</sup> مخفف کتابخانه ماتریکس<sup>۳</sup> است و یکی از پرکاربردترین زبان‌های برنامه‌نویسی شناخته شده است [۱۱۹]. متلب امکان محاسبات فنی و عددی، رسم داده‌ها، پیاده‌سازی الگوریتم‌ها، دست‌کاری ماتریس، یادگیری ماشین، شبیه‌سازی، پردازش تصویر و غیره را فراهم می‌سازد. در واقع یک محیط آسان برای ادغام برنامه‌نویسی، ترسیم و محاسبه است. متلب دارای کتابخانه‌های کاربردی فراوانی است که یادگیری و اجرای روش‌های متفاوت را برای کاربر آسان می‌کند. در این رساله به‌منظور پیش‌پردازش داده‌ها، مدل‌سازی شبکه عصبی مصنوعی و ارزیابی مدل‌های توسعه یافته از نرم‌افزار متلب ۲۰۱۷a استفاده شد.

<sup>1</sup>User-friendly

<sup>2</sup>MATLAB

<sup>3</sup>Matrix Laboratory

## ۲-۱-۶ نرم افزار Origin Lab

Origin Lab نرم‌افزاری قدرتمند است که برای تجزیه و تحلیل داده‌های آماری، پردازش توابع ریاضی و ویژوال بیسیک و رسم نمودارهای گرافیکی توسعه یافته است [۱۲۰]. از کاربردهای این نرم‌افزار می‌توان به مواردی همچون تجزیه و تحلیل داده‌ها و رسم نمودارهای آماری، قابلیت به اشتراک‌گذاری داده‌ها، نتایج و نمودارها با سایر نرم‌افزارها همچون اکسل، رسم نمودارهای دو بعدی و سه بعدی از توابع و داده‌های آماری، ذخیره نمودار با کیفیت بالا اشاره کرد. در این رساله برای نمایش برخی از نمودارها از نرم‌افزار Origin Lab استفاده شد.

## ۲-۱-۷ نرم‌افزار اتوداک<sup>۲</sup>

Autodock4.2 یک نرم‌افزار رایگان شبیه‌سازی و مدل‌سازی مولکولی است. این نرم‌افزار بر روی سیستم عامل به‌آسانی قابل اجرا می‌باشد. برنامه Autodock4.2 یک روش خودکار توسعه یافته برای پیش‌بینی برهم‌کنش<sup>۳</sup> لیگاندها با گیرنده‌های بیوماکرومولکولی است. هدف اصلی استفاده از این نرم‌افزار طراحی ترکیبات فعال زیستی در زمینه طراحی دارو به کمک رایانه می‌باشد که در آن با بررسی اتصال گیرنده - لیگاند، بهترین موقعیت لیگاند در جایگاه فعال گیرنده به دست می‌آید [۱۲۱].

## ۲-۱-۸ نرم افزار ویور لایت<sup>۴</sup>

ViewerLite یک برنامه رایگان برای مشاهده، ویرایش، بررسی وجود پیوندهای هیدروژنی و هیدروفوبی گیرنده- لیگاند، تخمین طول پیوند و زوایای پیوند مولکول‌های شیمیایی و همچنین تجزیه و تحلیل ساختارهای کریستالوگرافی پروتئین‌ها است [۱۲۲]. در این رساله، از نرم‌افزار ViewerLite نسخه ۵،

---

<sup>1</sup>Excel

<sup>2</sup>AutoDock

<sup>3</sup>Interaction

<sup>4</sup>ViewerLite

برای آماده‌سازی پروتئین و لیگاند و تعیین مختصات جایگاه فعال پروتئین قبل از اجرای فرایند داکینگ مولکولی استفاده شد.

## ۲-۱-۹ نرم افزار بیوویا دی اس<sup>۱</sup>

BIOVIA Discovery Studio برای مدل‌سازی و شبیه‌سازی در علوم دارویی بسیار کارآمد است. این نرم‌افزار برای مطالعه پروتئین، پپتید، لیگاندهای کوچک، لیپیدها و کربوهیدرات‌ها ایجاد شده است. نرم‌افزار مدل‌سازی و شبیه‌سازی BIOVIA Discovery Studio برای محققین امکان شبیه‌سازی، تجسم و تجزیه و تحلیل سیستم‌های شیمیایی و بیولوژیکی و آنالیز نتایج داکینگ مولکولی را فراهم می‌سازد. در این رساله برای نمایش برهم‌کنش‌های گیرنده-لیگاند در جایگاه فعال پروتئین از نرم افزار بیوویا دی اس استفاده شده است [۱۲۳].

## ۲-۱-۱۰ نرم افزار VMD<sup>۲</sup>

VMD محیط گرافیکی مناسبی برای شبیه‌سازی، نمایش و تجزیه و تحلیل سیستم‌های بیولوژیکی است. این نرم‌افزار قادر به فراخوانی فایل‌هایی با پسوند PDB می‌باشد و با استفاده از این نرم‌افزار می‌توان کمپلکس گیرنده-لیگاند را مورد بررسی قرار داد. به طوری که برهم‌کنش‌های متفاوت هیدروفیلی و هیدروفوبی با رنگ‌آمیزی‌های مختلف قابل مشاهده است.

---

<sup>1</sup>BIOVIA Discovery Studio

<sup>2</sup>Visual Molecular Dynamics

## ۲-۲ پیش‌بینی فعالیت دارویی برخی از مشتقات استانیلید / استامید

### به‌عنوان بازدارنده‌های ایدز با استفاده از مدل SCAD-ANN

#### ۲-۲-۱ مقدمه

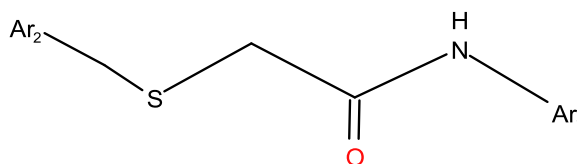
سندرم نقص ایمنی اکتسابی (ایدز) یک بیماری مزمن و تهدید کننده برای زندگی بشریت است که توسط ویروس نقص ایمنی انسانی (HIV) ایجاد می‌شود. HIV با آسیب رساندن به سیستم ایمنی بدن، توانایی بدن را برای مبارزه با عفونت و بیماری مختل می‌کند. HIV یک عفونت مقاربتی است و همچنین می‌تواند از طریق تماس با خون آلوده یا از مادر به کودک در دوران بارداری، زایمان یا شیردهی سرایت کند. هیچ درمان قطعی برای ایدز وجود ندارد، اما داروها می‌توانند به‌طور چشم‌گیری پیشرفت بیماری را کاهش دهند. این داروها مرگ و میر ناشی از ایدز را در بسیاری از کشورهای توسعه یافته کاهش داده است [۱۲۴]. ترکیبات فعال دارویی متفاوتی به‌عنوان بازدارنده‌های ایدز، مورد استفاده محققین دارویی قرار گرفته‌اند که از بین بازدارنده‌های بالقوه، NNRTIs دارای برخی ویژگی‌های مطلوب مانند سمیت کم، فعالیت بالا و انتخاب پذیری مناسب هستند [۱۲۵، ۱۲۶]. مشتقات Arylazolythioacetamide/acetanilide به‌عنوان گروه جدیدی از NNRTI ها، با خواص ساختاری مناسب، فعالیت ضد ایدز ویژه‌ای را در برابر سویه‌های جهش یافته نشان می‌دهند. در نتیجه طراحی و سنتز این ترکیبات حائز اهمیت بوده و مورد توجه محققین قرار گرفته است. این بازدارنده‌ها معمولاً از طریق اصلاح ساختاری پنج ترکیب رهبر طراحی و سنتز می‌شوند [۱۲۷-۱۳۳]. از آنجایی که سنتز و ارزیابی تجربی فعالیت دارویی ترکیبات فعال زمان‌بر و پرهزینه است، توسعه روش‌های تئوری از قبیل ساخت مدل QSAR برای ارزیابی و پیش‌بینی فعالیت‌های دارویی ترکیبات مشابه بسیار مهم است و باعث صرفه‌جویی در زمان، هزینه و به‌کارگیری نیروی انسانی می‌شود. با توجه به

<sup>1</sup>Sexually transmitted infection (STI)

اهمیت موضوع، در این بخش مدل QSAR برای پیش‌بینی فعالیت ترکیبات مورد مطالعه و پیشنهاد ترکیبات مشابه جدید بالقوه به کمک داکینگ مولکولی توسعه یافت. از این‌رو مدل شبکه عصبی جفت شده با روش انتخاب متغیر انقباضی SCAD برای پیش‌بینی فعالیت دارویی ( $pEC_{50}$ ) برخی از مشتقات استانیلید/استامید به‌عنوان بازدارنده‌های ایدز ساخته شد. مدل توسعه یافته با توصیف‌کننده‌های منتخب روش SCAD کمک شایانی را در راستای پیشنهاد ترکیبات فعال مشابه جدید نموده است. در ادامه این بخش با جزئیات بیشتری به مراحل ساخت مدل SCAD-ANN برای این دسته از بازدارنده‌ها پرداخته شده است.

## ۲-۲-۲ مجموعه داده‌ها

در این تحقیق مجموعه‌ای ۵۷ تایی از بازدارنده‌های ایدز از مجموعه مشتقات Arylazolythioacetamide/acetanilide جمع‌آوری شد و برای مدل‌سازی QSAR مورد بررسی قرار گرفت [۱۳۴-۱۳۶]. ساختار پایه ترکیبات مورد مطالعه در شکل ۱-۲ نشان داده شده است. ساختارهای شیمیایی ترکیبات و فعالیت‌های دارویی مربوطه در مقیاس لگاریتمی ( $pEC_{50} = \log \frac{1}{EC_{50}}$ ) محاسبه و در جدول ۱-۲ خلاصه شده است.  $EC_{50}$  غلظت مؤثر (بر حسب مولار) مورد نیاز برای مهار تکثیر HIV است [۱۳۷]. در این مطالعه از  $pEC_{50}$  با محدوده از ۴/۶۲ تا ۷/۷۴ به‌عنوان متغیر وابسته در ساخت مدل QSAR استفاده شده است. تقسیم‌بندی مجموعه داده‌ها با استفاده از الگوریتم KS انجام شد و داده‌ها به سه دسته آموزش (۳۵ ترکیب)، ارزیابی (۱۱ ترکیب) و آزمون (۱۱ ترکیب) تقسیم شدند و در مراحل مختلف مدل‌سازی QSAR مورد استفاده قرار گرفتند.



شکل ۱-۲ ساختار پایه ترکیبات Arylazolythioacetamide/acetanilide مورد مطالعه

جدول ۱-۲ ساختار ترکیبات شیمیایی به همراه مقادیر واقعی و پیش‌بینی شده pEC<sub>50</sub>

ردیف	Ar <sub>1</sub>	Ar <sub>2</sub>	pE <sub>50</sub> واقعی	pE <sub>50</sub> پیش‌بینی شده
۱	2-fluorophenyl	2-(3-(2-Chlorophenyl) pyrazin-2-yl)	۵/۲۸	۵/۴
۲ <sup>t</sup>	2-Chlorophenyl	2-(3-(2-chlorophenyl) pyrazin-2-yl)	۵/۳۴	۵/۳۷
۳ <sup>v</sup>	2,4-dichlorophenyl	2-(3-(2-Chlorophenyl) pyrazin-2-yl)	۵/۲۸	۵/۱۶
۴	2-Bromophenyl	2-(3-(2-chlorophenyl) pyrazin-2-yl)	۵/۳۸	۵/۳۴
۵	2-Bromo-4-methylphenyl	2-(3-(2-chlorophenyl) pyrazin-2-yl)	۵/۳۷	۵/۲۹
۶ <sup>t</sup>	2-Bromo-4-chlorophenyl	2-(3-(2-chlorophenyl) pyrazin-2-yl)	۵/۲۹	۵/۱۷
۷ <sup>v</sup>	4-Acetyl-2-bromophenyl	2-(3-(2-chlorophenyl) pyrazin-2-yl)	۵/۶	۵/۳
۸ <sup>t</sup>	Methyl 3-bromo-4-benzoat	(2-(3-(2-chlorophenyl) pyrazin-2-yl)	۵/۴۳	۵/۱۴
۹ <sup>v</sup>	2-nitrophenyl	2-(3-(2-Chlorophenyl) pyrazin-2-yl)	۵/۰۵	۵/۲
۱۰	4-methyl-2-nitrophenyl	2-(3-(2-Chlorophenyl) pyrazin-2-yl)	۵/۷۷	۵/۷۴
۱۱ <sup>t</sup>	2-chloropyridin-3-yl	2-(3-(2-Chlorophenyl) pyrazin-2-yl)	۵/۵۴	۵/۳۵
۱۲	2-chloropyridin-3-yl (3,4-	2-(3-(2-Chlorophenyl) pyrazin-2-yl)	۵/۶۲	۵/۷
۱۳	dihydroisoquinolin-2(1H)-yl) ethenone	2-(3-(2-Chlorophenyl) pyrazin-2-yl)	۵/۳۹	۵/۲۸
۱۴ <sup>v</sup>	Phenyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۶/۳۸	۶/۲۷
۱۵	2-chloro Phenyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۶/۰۹	۶/۵۵
۱۶	2-fluoro Phenyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۶/۶۸	۶/۳۷
۱۷	2-bromo Phenyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۷/۰۴	۶/۹۶
۱۸ <sup>v</sup>	2-bromo 4-methylPhenyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۶/۷۵	۶/۶۱
۱۹ <sup>t</sup>	4-methyl-2-nitrophenyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۷/۳۲	۷/۶۹
۲۰	2-nitrophenyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۷/۶۴	۷/۵۸
۲۱	4-acetyl-2-bromophenyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۷/۰۷	۶/۸۸
۲۲	2-chloropyridin-3-yl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۷/۱۳	۷/۱۹
۲۳ <sup>v</sup>	2,4-dichlorophenyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۶/۵۹	۶/۵۸
۲۴ <sup>t</sup>	o-tolyl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۶/۵۳	۶/۱۳
۲۵	ethyl 3-bromo-4-benzoate	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۶/۷۲	۶/۷۱
۲۶	3-bromo-5-methylpyridin-2-yl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۶/۶۳	۶/۷۸
۲۷	pyridin-2-yl	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۵/۴۸	۶/۰۲
۲۸ <sup>v</sup>	Methyl-3-thiophene-2-carboxylate	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۵/۵۴	۵/۴۱
۲۹ <sup>t</sup>	1-(3,4-dichlorophenyl)ethan-1-one	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۵/۲۱	۴/۵۸
۳۰	1-(3,4-dihydroisoquinolin-2(1H)-yl)ethan-1-one	2-(5-(2-chlorophenyl)-1,2,4-triazin-6-yl)	۶/۱۴	۵/۸۵

## ادامه جدول ۱-۲

ردیف	Ar <sub>1</sub>	Ar <sub>2</sub>	pE <sub>50</sub> واقعی	pE <sub>50</sub> پیش‌بینی شده
۳۱	phenyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۵/۸۲	۵/۸۹
۳۲	2-chlorophenyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۶/۸۷	۶/۶۶
۳۳	2-fluoro Phenyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۶/۴۴	۶/۹
۳۴ <sup>۷</sup>	2-bromo Phenyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۶/۷۶	۶/۹۶
۳۵	2-bromo 4-methylPhenyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۶/۷۱	۶/۶۴
۳۶	4-methyl-2-nitrophenyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۷/۳۴	۷/۲۲
۳۷	2-nitrophenyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۷/۳۷	۷/۳۷
۳۸ <sup>۴</sup>	4-acetyl-2-bromophenyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۷/۲۹	۶/۷۸
۳۹	2-chloropyridin-3-yl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۷/۲۳	۶/۷۹
۴۰-۷	o-tolyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۶/۰۸	۶/۰۲
۴۱ <sup>۴</sup>	methyl 3-bromo-4-benzoate	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۷/۱۶	۶/۸۶
۴۲ <sup>۷</sup>	ethyl 3-bromo-4-benzoate	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۶/۶۸	۶/۹۷
۴۳	ethyl 3-chloro-4-benzoate	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۶/۶۹	۶/۷۸
۴۴	2-bromo-4-sulfamoylphenyl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۷/۷۴	۷/۷۶
۴۵ <sup>۴</sup>	naphthalen-1-yl	2-(5-(naphthalene-1-yl)-1,2,4-triazin-6-yl)	۶/۱۲	۵/۷۸
۴۶	2-bromo-4-sulfamoylphenyl	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۶/۷۴	۶/۷۷
۴۷ <sup>۷</sup>	Methyl 3-bromo-benzoate	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۵/۶۲	۵/۳۴
۴۸	Ethyl 3-bromo-benzoate	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۵/۳۵	۵/۳۶
۴۹ <sup>۴</sup>	Methyl 3-chloro-benzoate	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۵/۵۵	۵/۳
۵۰	Ethyl 3-chloro-benzoate	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۵/۳۴	۵/۲۳
۵۱	3-Chloro-benzoic acid	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۵/۰۹	۵/۲۴
۵۲	2-nitrophenyl	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۶/۰۴	۶/۰۴
۵۳	3-Chloro-N-methoxy-benzamide	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۶/۱۲	۶/۱
۵۴	3-Chloro-N-hydroxy-benzamide	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۶/۲	۶/۰۲
۵۵	Ethyl 2-(3-chloro-benzamide)acetate	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۶/۸۲	۶/۵۱
۵۶	6-Chloro-5-picolinic acid	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۴/۶۲	۴/۶۱
۵۷	Ethyl 2-(3-chloro-benzamide)propanoate	2-(5-(Naphthalen-1-yl)pyrimidin-4-yl)	۵/۵۶	۵/۸۶

\* ۷ و t به ترتیب نمایانگر داده‌های مجموعه ارزیابی و آزمون هستند و سایر ترکیبات در مجموعه آموزش قرار دارند.



## ۲-۲-۳ رسم و بهینه‌سازی ساختار Arylazolythioacetamide/acetanilide ها

ساختار شیمیایی مشتقات مورد مطالعه با استفاده از نرم‌افزار هایپرکم رسم شد و بهینه‌سازی ساختارها مطابق روش کار بخش (۱-۵-۳) تا رسیدن به حداقل انرژی بهینه شدند.

## ۲-۲-۴ استخراج توصیف‌کننده‌ها

ساختارهای بهینه ترکیبات مورد مطالعه در نرم‌افزار دراگون فراخوانی شدند و سپس برای هر ترکیب به تعداد ۳۲۲۴ توصیف‌کننده در ۲۲ دسته متفاوت محاسبه شدند.

## ۲-۲-۵ پیش‌پردازش و انتخاب توصیف‌کننده‌های مؤثر

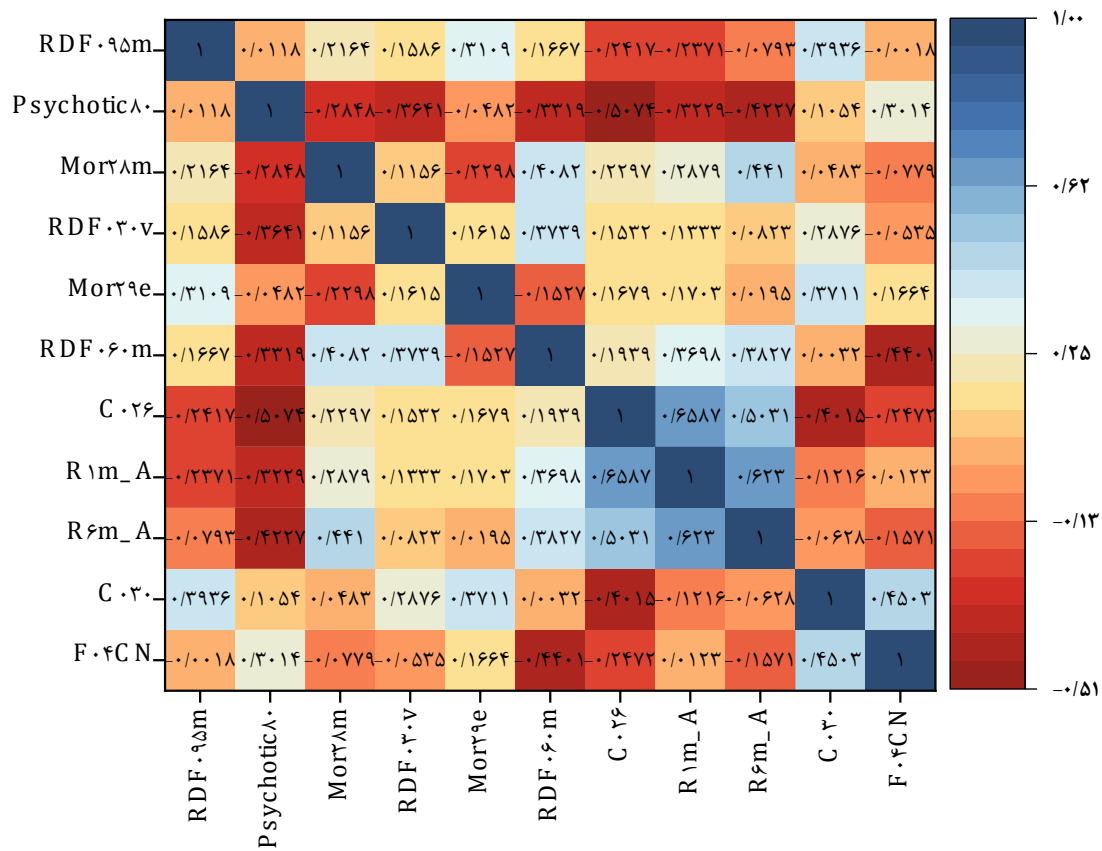
با توجه به تعداد زیاد توصیف‌کننده‌های محاسبه شده و احتمال وجود توصیف‌کننده‌های اضافی و فاقد اطلاعات مفید، فرآیند پیش‌پردازش و غربالگری توصیف‌کننده‌ها به‌منظور کاهش تعداد توصیف‌کننده‌ها و در نتیجه بهبود صحت پیش‌بینی، کاهش زمان مدل‌سازی و افزایش تفسیرپذیری مدل بر روی کل توصیف‌کننده‌های محاسبه شده اعمال شد. در این فرآیند، توصیف‌کننده‌هایی با مقادیر ثابت و تقریباً ثابت (توصیف‌کننده‌هایی با واریانس کم‌تر از ۰/۰۰۱) از مجموع توصیف‌کننده‌های محاسبه شده حذف شدند و در نهایت از بین دو توصیف‌کننده همبسته ( $R^2 > 0/90$ ) که دارای اطلاعات تقریباً یکسانی هستند، توصیف‌کننده با بیش‌ترین همبستگی با پاسخ، نگه داشته شد و توصیف‌کننده بعدی حذف شد. پس از کاهش اولیه توصیف‌کننده‌ها، فرآیند انتخاب متغیر انقباضی با روش ارزیابی تقاطعی ده فولد<sup>۱</sup> موجود در بسته ncvreg در نرم‌افزار R، روی داده‌های مجموعه آموزش و ارزیابی اجرا شد تا مؤثرترین توصیف‌کننده‌ها انتخاب شود. با اجرای روش SCAD، تعداد ۱۱ توصیف‌کننده مربوط به  $\lambda_{\min}$  ( $\lambda$  دارای کم‌ترین خطای ارزیابی تقاطعی) انتخاب شدند. نام و نوع توصیف‌کننده‌های منتخب روش SCAD در جدول ۲-۲ خلاصه شده‌اند. به‌منظور

<sup>۱</sup>Ten fold cross validation (10-fold-CV-SCAD)

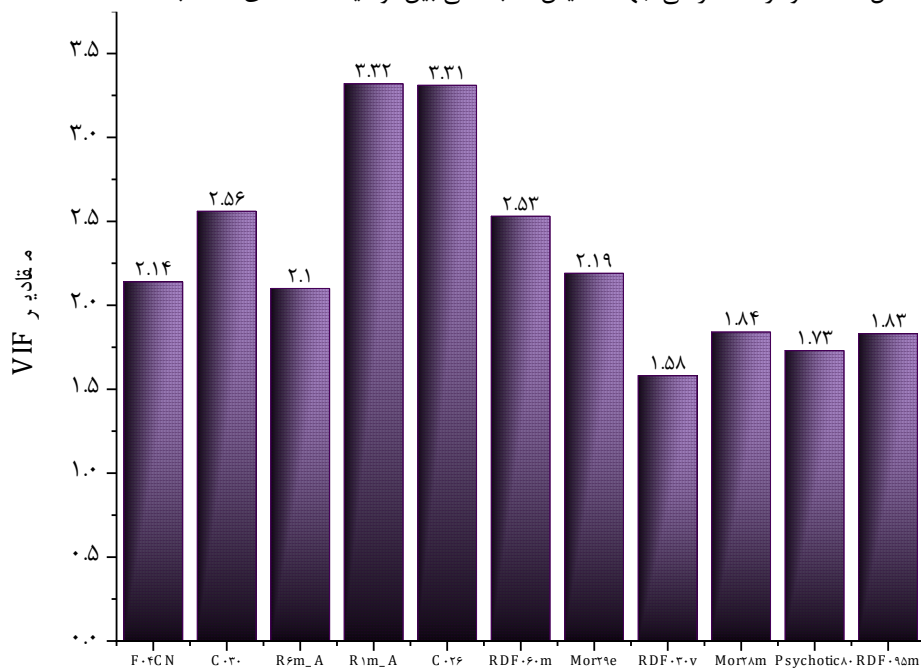
بررسی دقیق‌تر توصیف‌کننده‌های منتخب با روش SCAD، احتمال وجود همبستگی و هم‌خطی بین توصیف‌کننده‌ها، به ترتیب با محاسبه مقادیر ضریب همبستگی بین دو توصیف‌کننده و مقادیر افزایش تورم واریانس (VIF) توصیف‌کننده (مطابق رابطه ۱-۱۰) مطالعه گردید. به این منظور نمودارهای نقشه رنگی و VIF برای توصیف‌کننده‌های منتخب رسم شد و نتایج حاصل در شکل ۲-۲ و شکل ۳-۲ نمایش داده شده‌اند. نتایج مربوط به نمودار نقشه رنگی (شکل ۲-۲) نشان می‌دهد که همبستگی معناداری بین توصیف‌کننده‌های انتخاب شده با روش SCAD وجود ندارد. علاوه بر این نمودار VIF (شکل ۳-۲) نیز نشان می‌دهد که همه توصیف‌کننده‌های منتخب با روش SCAD مقادیر VIF کم‌تر از ۱۰ دارند، این شواهد بیانگر این است که بین توصیف‌کننده‌های منتخب با روش SCAD هم‌خطی قابل ملاحظه‌ای وجود ندارد [۱۳۸، ۱۳۹].

جدول ۲-۲ توصیف‌کننده‌های منتخب SCAD

ردیف	نماد	طبقه‌بندی	معنا	اهمیت
۱	F04[C-N]	2D frequency fingerprints	Frequency of C-N at topological distance 04	۰/۳۲۴
۲	C030	Atom centered fragments	X—CH—X	۰/۲۵۶
۳	R6m_A	GETAWAY	R autocorrelation of lag 6 / weighted by atomic masses	۰/۲۰۲
۴	R1m_A	GETAWAY	R autocorrelation of lag 1 / weighted by atomic masses	۰/۱۷
۵	C026	Atom centered fragments	R—CX—R	۰/۱۲۳
۶	RDF060m	RDF	Radial distribution function -6.0 / weighted by atomic masses	۰/۱۰
۷	Mor29e	3D-MoRSE	3D- MoRSE- signal 29/ weighted by atomic van der waals volumes	۰/۰۹۳
۸	RDF030v	RDF	Radial distribution function -3.0/ weighted by atomic masses	۰/۰۸۹
۹	Mor28m	3D-MORSE	3D- MoRSE- signal 29/ weighted by atomic masses	۰/۰۸۷
۱۰	Psychotic80	Molecular properties	Ghose- viswandhan- Wendoloski antipsychotic – like index at 80%	۰/۰۸۶
۱۱	RDF095m	RDF	Radial distribution function -9.5/ weighted by atomic masses	۰/۰۸۴



شکل ۲-۲ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب SCAD



توصیف‌کننده‌ها

شکل ۲-۳ نمودار مقادیر VIF توصیف‌کننده‌های منتخب SCAD

## ۲-۲-۶ مدل سازی شبکه عصبی با استفاده از توصیف کننده های منتخب SCAD

به منظور ساخت مدل QSAR غیرخطی از روش شبکه عصبی مصنوعی استفاده شد. در این مطالعه، یک مدل شبکه عصبی مصنوعی (ANN) پیشخور با الگوریتم آموزشی پس انتشار خطا مورد استفاده قرار گرفت. برنامه آن در یک m-file در نرم افزار MATLAB نوشته شد. با توجه به توصیه مقالات منتشر شده، استفاده از یک لایه پنهان در مدل ANN، در مطالعات شیمی محاسباتی کفایت می کند [۱۴۰، ۱۴۱]. بنابراین، مدل های ANN سه لایه ای شامل یک لایه ورودی، یک لایه پنهان و یک لایه خروجی برای بهینه سازی پارامترهای ANN و نهایتاً ساخت مدل نهایی استفاده شد. برای به دست آوردن بهترین مدل ANN، تمامی پارامترهای مؤثر بر عملکرد پیش بینی مدل، مانند تعداد ورودی ها، گره های لایه پنهان، دور آموزش و توابع آموزش و انتقال بهینه شدند. برای این منظور، چهار مدل شبکه عصبی مصنوعی مختلف با استفاده از دو الگوریتم آموزشی متفاوت لونیگ - مارکوارت (LM) (تابع آموزشی trainlm در جعبه ابزار متلب) و تنظیم بایزین (BR) (تابع آموزشی trainbr در جعبه ابزار متلب) و دو تابع انتقال لگاریتم سیگموئیدی و تانژانت هایپربولیک سیگموئیدی (در جعبه ابزار برنامه متلب به ترتیب با توابع logsig و tansig شناخته می شوند) طراحی شدند. در تمام مدل های ANN طراحی شده، تابع خطی (purlin) به عنوان تابع انتقال خروجی استفاده شد. سایر پارامترهای مدل های ANN، مانند تعداد ورودی، تعداد گره ها و تعداد دوره های آموزش به طور همزمان بهینه شدند. در بهینه سازی تعداد ورودی های ANN، زیرمجموعه هایی با تعداد توصیف کننده در بازه ۲ تا ۱۱ از میان توصیف کننده های انتخاب شده با SCAD به عنوان ورودی ANN در نظر گرفته شدند. با توجه به تعداد زیاد این زیرمجموعه ها (حدود  $10^{10}$  زیرمجموعه دوتایی تا یازده تایی)، طراحی و استفاده از تمام توصیف کننده های تصادفی ایجاد شده به عنوان ورودی ANN عملاً غیرممکن است. بنابراین طبق توضیحات مندرج در بخش ۱-۵-۷-۲ از چیدمان شبکه عصبی برای تعیین اهمیت توصیف کننده های منتخب روش SCAD به عنوان ورودی شبکه عصبی مورد استفاده قرار گرفت.

به طوری که ابتدا شبکه عصبی با ۱۱ زیرمجموعه به عنوان ورودی شبکه عصبی مصنوعی تعریف شد. برای این منظور، ابتدا مدل های ANN (۴ معماری ANN با توجه به دو تابع آموزش و انتقال متفاوت) با کل توصیف کننده های منتخب روش SCAD بهینه شدند. مدل SCAD-LM-ANN با معماری ۱-۴-۱۱ و با دور آموزشی ۵، تابع آموزش LM و تابع انتقال tansig دارای حداقل خطای مجموعه ارزیابی است. پس از شناسایی مدل بهینه شبکه عصبی توسعه یافته با ۱۱ توصیف کننده، از مدل بهینه برای تعیین اهمیت توصیف کننده های منتخب روش SCAD استفاده شد. برای تخمین اهمیت توصیف کننده  $i$  ام، همه ۱۱ توصیف کننده وارد مدل ANN بهینه شدند، در حالی که مقادیر توصیف کننده  $i$  ام با مقادیر تصادفی در محدوده تغییرات مقادیر واقعی توصیف کننده  $i$  جایگزین شدند. مدل ANN با ماتریس ورودی دست کاری شده آموزش داده شد و مقادیر  $PEC_{50}$  داده های مجموعه ارزیابی پیش بینی شد و مقدار MSE مجموعه ارزیابی محاسبه شد. این فرآیند برای همه توصیف کننده ها تکرار شد، به طوری که مقادیر همه توصیف کننده ها یکبار با مقادیر تصادفی جایگزین شدند و MSE مجموعه ارزیابی محاسبه شد. در نتیجه، ۱۱ MSE برای ۱۱ توصیف کننده به دست آمد. بالاترین MSE نشان می دهد که مدل ANN در غیاب آن توصیف کننده، خطای بیش تری را در مدل بهینه متحمل می شود و چنین توصیف کننده ای بیشترین اهمیت را در توسعه مدل ANN دارد. توصیف کننده ها بر اساس اهمیت آن ها (مقادیر MSE از زیاد به کم) در جدول ۲-۲ مرتب شده اند. از مقادیر اهمیت محاسبه شده ۱۱ توصیف کننده، برای طراحی ۱۰ زیرمجموعه استفاده شد. به طوری که زیرمجموعه اول شامل توصیف کننده های مهم اول و دوم بود و زیرمجموعه های بعدی به ترتیب با افزودن یک توصیف کننده دیگر بر اساس ترتیب اهمیت توصیف کننده ها به زیرمجموعه های قبلی ایجاد شدند. پس از ایجاد زیر مجموعه های ورودی با استفاده از توصیف کننده های منتخب روش SCAD، تعداد ورودی ها، تعداد نورون های لایه پنهان و دور آموزشی به طور هم زمان بهینه شدند. در بهینه سازی پارامترهای ANN، تعداد ورودی ها با استفاده از زیرمجموعه های ایجاد شده با ۲ تا ۱۱ توصیف

کننده به عنوان ورودی‌های ANN در بازه ۲ تا ۱۱ تغییر یافت. علاوه بر این، تعداد نوروها در لایه پنهان و دور آموزشی به ترتیب در محدوده ۲ تا ۱۰ (با گام ۱) و ۵ تا ۵۰ (با گام ۵) به طور همزمان تغییر یافت. تمام مدل‌های ممکن و ایجاد شده ANN، با استفاده از مجموعه داده‌های آموزش (شامل ۳۵ ترکیب) آموزش داده شدند و برای پیش‌بینی مقادیر  $pEC_{50}$  مربوط به ۱۱ ترکیب موجود در مجموعه ارزیابی استفاده شدند. مقادیر MSE مجموعه ارزیابی برای تمام مدل‌های آموزش دیده محاسبه و به عنوان معیاری برای انتخاب بهترین مدل ANN مورد استفاده قرار گرفت. با توجه به نتایج، معماری شبکه و مقادیر MSE برای چهار مدل ANN با توابع آموزش و انتقال متفاوت که دارای کمترین مقادیر MSE بودند در جدول ۲-۳ خلاصه شده‌اند. نتایج نشان می‌دهند که مدل ANN با ۵ توصیف کننده پر اهمیت و منتخب روش SCAD به عنوان ورودی، ۵ گره در لایه پنهان و ۵ دور آموزش با تابع انتقال لگاریتم سیگموئیدی و الگوریتم آموزشی LM دارای کمترین مقدار MSE برای پیش‌بینی داده‌های مجموعه ارزیابی می‌باشد. بنابراین این مدل با نماد LM-SCAD-ANN نشان داده شد و با معماری ۱-۵-۵ به عنوان مدل برتر برای پیش‌بینی مقادیر  $pEC_{50}$  ترکیبات مورد مطالعه انتخاب شد.

جدول ۲-۳ ساختارهای شبکه‌های توسعه یافته با توصیف کننده‌های منتخب SCAD با کمترین MSE مجموعه ارزیابی

تعداد توصیف کننده	تابع آموزش	تابع انتقال	تعداد گره	تعداد دور آموزش	MSE <sub>validation</sub>	R <sup>2</sup> <sub>validation</sub>
۱۰	تنظیم یازین	لگاریتم-سیگموئید	۲	۱۵	۰/۰۶	۰/۸۹
۵	لونیبرگ-مارکوارت	لگاریتم-سیگموئید	۵	۵	۰/۰۳	۰/۹۷
۶	تنظیم یازین	تانژانت-سیگموئید	۴	۵	۰/۰۶	۰/۸۹
۹	لونیبرگ-مارکوارت	تانژانت-سیگموئید	۳	۲۰	۰/۰۴	۰/۹۱

برای مقایسه عملکرد روش SCAD در انتخاب مؤثرترین توصیف کننده‌ها، از روش متداول انتخاب متغیر رگرسیون گام به گام (SR) استفاده شد. SR بر روی مجموعه داده‌های ارزیابی و آموزش اعمال شد. تعداد ۱۱ توصیف کننده شامل (F04CN, R1m\_A, RDF095m, HATS5v, F09CN, C030, C005,

Mor12v, Mor22u, RDF080m, DISPv انتخاب شدند. زیرمجموعه‌هایی شامل ۲ تا ۱۱ توصیف کننده بر اساس اهمیت آن‌ها در مدل ANN (مطابق با بخش ۱-۵-۷-۲) ایجاد شدند و به‌عنوان ورودی برای طراحی و بهینه سازی مدل‌های شبکه عصبی استفاده شدند. نتایج نشان داد مدل ANN با تابع آموزش LM و با استفاده از زیر مجموعه‌های ۱۱ تایی از توصیف کننده‌های SR با تعداد ۲ گره در لایه پنهان و ۵ دور آموزش حداقل مقدار MSE برای مجموعه ارزیابی ایجاد نمود. مدل بهینه SR-LM-ANN برای پیش‌بینی داده‌های مجموعه آزمون ( $pEC_{50}$ ) استفاده شد.

## ۲-۲-۷ ارزیابی مدل SCAD-LM-ANN

یکی از اساسی‌ترین مراحل مدل‌سازی، بررسی قدرت پیش‌بینی مدل با استفاده از تکنیک‌های مختلف آماری است. در این تحقیق، قدرت پیش‌بینی، اعتبار و تعمیم پذیری مدل بهینه توسعه یافته SCAD-LM-ANN با استفاده از پیش‌بینی داده‌های مجموعه آزمون، پیش‌بینی پاسخ کل ترکیبات با تکنیک LOO، نمودار باقی‌مانده‌ها، محاسبه پارامترهای آماری متفاوت، آنالیز دامنه کاربرد و آزمون  $Y$ -تصادفی مورد ارزیابی قرار گرفت.

## ۲-۲-۷-۱ ارزیابی مدل SCAD-LM-ANN با استفاده از پیش‌بینی داده‌های مجموعه آزمون

اعتبار و اهمیت مدل با استفاده از داده‌های مجموعه آزمون که در حین انتخاب متغیر و مدل‌سازی حضور نداشتند، سنجیده شد. به‌این منظور فعالیت دارویی داده‌های مجموعه آزمون با استفاده از مدل توسعه یافته بهینه SCAD-LM-ANN آموزش دیده در شرایط بهینه و با معماری ۱-۵-۵ پیش‌بینی شد. مقادیر پیش‌بینی شده فعالیت ترکیبات موجود در مجموعه آزمون در جدول ۲-۴ آورده شده است. مقادیر خطای کم اغلب ترکیبات موجود در مجموعه آزمون، نشان‌دهنده قدرت پیش‌بینی مناسب مدل توسعه یافته است. برای بررسی بیشتر قدرت پیش‌بینی مدل در مجموعه آزمون، مقادیر  $pEC_{50}$  پیش‌بینی شده به‌وسیله مدل بر حسب مقادیر واقعی آن‌ها رسم شد. شکل ۲-۴ نشان‌دهنده مقدار  $R^2$  قابل قبول مدل بهینه ANN است.

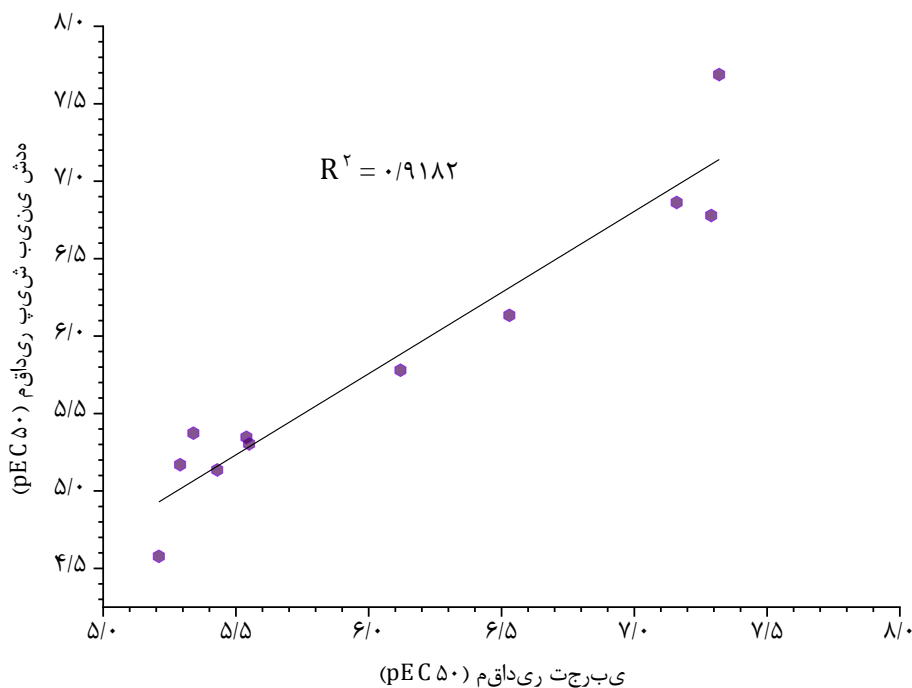
با توجه به این که  $R^2$  مربوط به مجموعه آموزش و ارزیابی برای مدل SCAD-LM-ANN به ترتیب برابر با ۰/۹۱ و ۰/۹۴ است بنابراین با مقایسه مقادیر مذکور با مقدار  $R^2$  مجموعه آزمون، قدرت پیش بینی و تعمیم پذیری مدل SCAD-LM-ANN اثبات می شود.

همان طور که در بخش ۲-۲-۶ گفته شد، به منظور مقایسه عملکرد SCAD-LM-ANN، از روش SR-LM-ANN استفاده شد. مقادیر MSE و  $R^2_{test}$  برای مجموعه آزمون به ترتیب برابر با ۰/۴۲ و ۰/۴۸ به دست آمد. نتایج نشان می دهد که قدرت پیش بینی مدل SR-ANN رضایت بخش نیست و SCAD-LM-ANN توانایی پیش بینی بسیار بهتری را نسبت به SR-LM-ANN دارد.

جدول ۴-۲ نتایج حاصل از ارزیابی مدل SCAD-ANN با استفاده از مجموعه آزمون

شماره ترکیب در مجموعه داده‌ها	pEC <sub>۵۰</sub>		درصد خطا
	مقدار واقعی	مقدار پیش بینی شده	
۲	۵/۳۴	۵/۳۷	۰/۶۱
۶	۵/۲۹	۵/۱۷	-۲/۲۸
۸	۵/۴۳	۵/۱۴	-۵/۴۲
۱۱	۵/۵۴	۵/۳۵	-۳/۴۹
۱۹	۷/۳۲	۷/۶۹	۵/۰۴
۲۴	۶/۵۳	۶/۱۳	-۶/۰۸
۲۹	۵/۲۱	۴/۵۸	-۱۲/۱۵
۳۸	۷/۲۹	۶/۷۸	-۷/۰۱
۴۱	۷/۱۶	۶/۸۶	-۴/۱۶
۴۵	۶/۱۲	۵/۷۸	-۵/۵۷
۴۹	۵/۵۵	۵/۳	-۴/۴۵





شکل ۲-۴ نمودار تغییرات مقادیر پیش‌بینی شده  $pEC_{50}$  به وسیله مدل SCAD-LM-ANN در شرایط بهینه در مقابل مقادیر تجربی برای داده‌های مجموعه آزمون

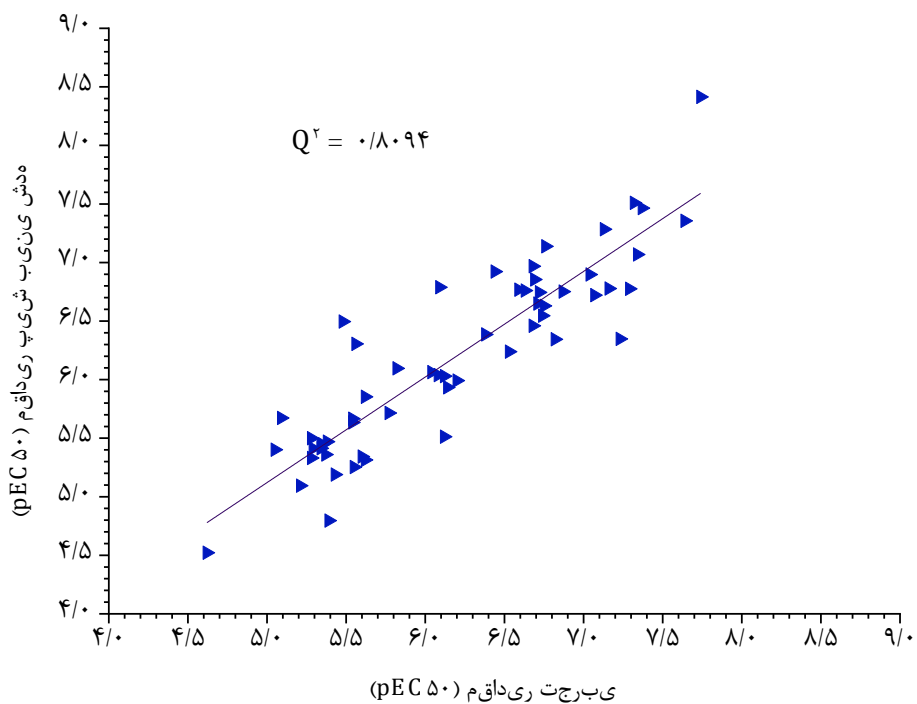
## ۲-۲-۷-۲ ارزیابی مدل SCAD-ANN با استفاده از روش رد مرحله‌ای تک تک

در این بخش، یکی از تکنیک‌های قدرتمند برای ارزیابی مدل بهینه و بررسی قدرت پیش‌بینی مدل به نام تکنیک رد مرحله‌ای تک تک (LOO) برای پیش‌بینی همه داده‌های مورد مطالعه به عنوان داده آزمون مورد استفاده قرار گرفت. در این مطالعه هر بار یکی از داده‌های مورد مطالعه کنار گذاشته شد و مدل SCAD-LM-ANN با سایر داده‌ها آموزش داده شد و فعالیت دارویی ترکیب کنار گذاشته با مدل آموزشی داده شده، پیش‌بینی گردید. این فرایند برای همه ترکیبات انجام شد و همه داده‌ها یک‌بار به عنوان داده آزمون در نظر گرفته شدند. از این رو فعالیت دارویی همه ترکیبات تکرار شد به گونه‌ای که تکنیک LOO توسط مدل SCAD-LM-ANN پیش‌بینی شدند و نتایج حاصله در جدول ۲-۵ آورده شده‌اند. برای مطالعه بیش‌تر قدرت پیش‌بینی مدل، نمودار مقادیر پیش‌بینی شده  $pEC_{50}$  تمام ترکیبات بر حسب مقادیر تجربی آن‌ها رسم شد. مقدار  $Q^2_{LOO}$  که در شکل ۲-۵ نشان داده شده است بیش‌تر از مقدار قابل قبول ( $Q^2 > 0.5$ )

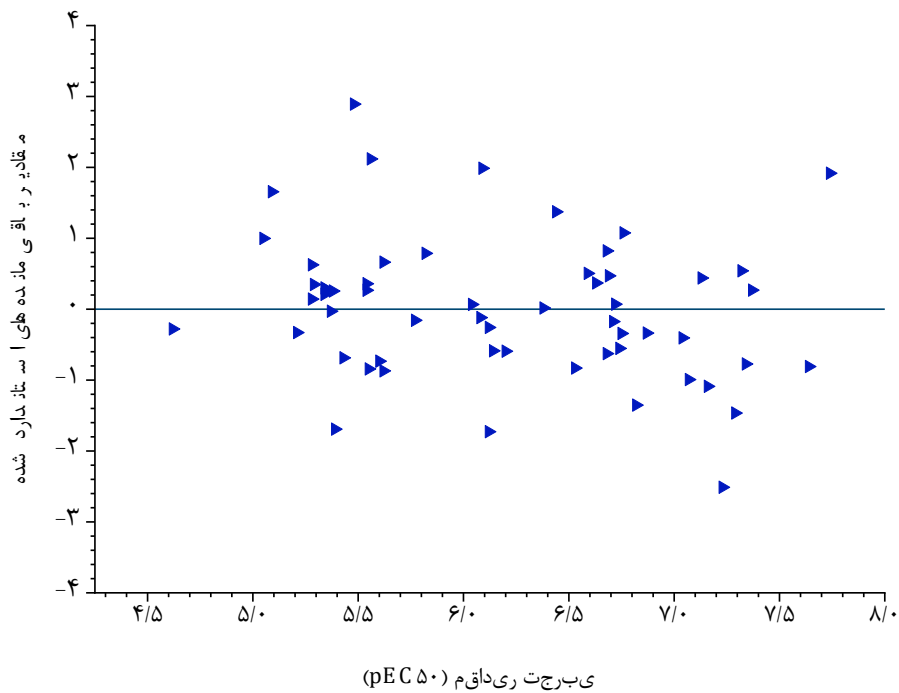
بوده که نشان‌دهنده این است که مدل توسعه یافته از پایداری مناسبی برخوردار است. علاوه بر این، به‌منظور بررسی عدم وجود خطای سیستماتیک در مدل توسعه یافته، مقادیر باقی‌مانده بر حسب مقادیر تجربی رسم شد. نتایج حاصله شکل ۲-۶ نشان‌دهنده توزیع یکنواخت و تصادفی داده‌ها حول محور صفر می‌باشد که بیانگر عدم وجود خطای سیستماتیک در مدل شبکه عصبی توسعه یافته با استفاده از روش SCAD به‌عنوان روش انتخاب متغیر (SCAD-LM-ANN) است.

جدول ۵-۲ نتایج حاصل از ارزیابی مدل SCAD-LM-ANN به روش رد مرحله‌ای تک تک برای کل داده‌ها

شماره ترکیب	pEC <sub>50</sub>			شماره ترکیب	pEC <sub>50</sub>		
	مقدار واقعی	مقدار پیش‌بینی شده	درصد خطا		مقدار واقعی	مقدار پیش‌بینی شده	درصد خطا
۱	۵/۲۸	۵/۵	۴/۱۶	۳۰	۶/۱۴	۵/۹۳	-۳/۳۶
۲	۵/۳۴	۵/۴۱	۱/۳۸	۳۱	۵/۸۲	۶/۱	۴/۷۵
۳	۵/۲۸	۵/۳۳	۰/۹۵	۳۲	۶/۸۷	۶/۷۵	-۱/۷۲
۴	۵/۳۸	۵/۴۷	۱/۶۷	۳۳	۶/۴۴	۶/۹۲	۷/۴۹
۵	۵/۳۷	۵/۳۶	-۰/۱۹	۳۴	۶/۷۶	۷/۱۴	۵/۵۹
۶	۵/۲۹	۵/۴۱	۲/۳۱	۳۵	۶/۷۱	۶/۶۵	-۰/۹۱
۷	۵/۶	۵/۳۴	-۴/۶۱	۳۶	۷/۳۴	۷/۰۷	-۳/۷
۸	۵/۴۳	۵/۱۹	-۴/۴۴	۳۷	۷/۳۷	۷/۴۶	۱/۲۸
۹	۵/۰۵	۵/۴	۶/۹۵	۳۸	۷/۲۹	۶/۷۸	-۷/۰۶
۱۰	۵/۷۷	۵/۷۱	-۰/۹۶	۳۹	۷/۲۳	۶/۳۵	-۱۲/۲۱
۱۱	۵/۵۴	۵/۶۷	۲/۲۸	۴۰	۶/۰۸	۶/۰۴	-۰/۶۸
۱۲	۵/۶۲	۵/۸۵	۴/۱۴	۴۱	۷/۱۶	۶/۷۸	-۵/۳۴
۱۳	۵/۳۹	۴/۸	-۱۱/۰۲	۴۲	۶/۶۸	۶/۹۷	۴/۳۲
۱۴	۶/۳۸	۶/۳۹	۰/۰۹	۴۳	۶/۶۹	۶/۸۶	۲/۴۷
۱۵	۶/۰۹	۶/۷۹	۱۱/۴۷	۴۴	۷/۷۴	۸/۴۱	۸/۷۱
۱۶	۶/۶۸	۶/۴۶	-۳/۲۹	۴۵	۶/۱۲	۵/۵۱	-۹/۹۲
۱۷	۷/۰۴	۶/۹	-۲/۰۲	۴۶	۶/۷۴	۶/۵۵	-۲/۸۸
۱۸	۶/۷۵	۶/۶۳	-۱/۷۸	۴۷	۵/۶۲	۵/۳۱	-۵/۴۴
۱۹	۷/۳۲	۷/۵۱	۲/۵۹	۴۸	۵/۳۵	۵/۴۵	۱/۸۲
۲۰	۷/۶۴	۷/۳۶	-۳/۷۲	۴۹	۵/۵۵	۵/۲۵	-۵/۳۴
۲۱	۷/۰۷	۶/۷۲	-۴/۹۴	۵۰	۵/۳۴	۵/۴۴	۱/۹۶
۲۲	۷/۱۳	۷/۲۸	۲/۱۷	۵۱	۵/۰۹	۵/۶۷	۱۱/۴۳
۲۳	۶/۵۹	۶/۷۷	۲/۶۹	۵۲	۶/۰۴	۶/۰۶	۰/۳۹
۲۴	۶/۵۳	۶/۲۴	-۴/۴۷	۵۳	۶/۱۲	۶/۰۳	-۱/۴۸
۲۵	۶/۷۲	۶/۷۴	۰/۳۶	۵۴	۶/۲	۵/۹۹	-۳/۳۶
۲۶	۶/۶۳	۶/۷۶	۱/۹۶	۵۵	۶/۸۲	۶/۳۴	-۶/۹۷
۲۷	۵/۴۸	۶/۵	۱۸/۵۴	۵۶	۴/۶۲	۴/۵۲	-۲/۱۳
۲۸	۵/۵۴	۵/۶۳	۱/۶۹	۵۷	۵/۵۶	۶/۳	۰/۷۴
۲۹	۵/۲۱	۵/۰۹	-۲/۲۲				



شکل ۲-۵ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌ها بر اساس تکنیک LOO در مقابل مقادیر تجربی



شکل ۲-۶ نمودار باقی‌مانده‌های حاصل از پیش‌بینی فعالیت دارویی ترکیبات با استفاده از تکنیک LOO و مقادیر تجربی برحسب مقادیر تجربی

## ۲-۲-۷-۳ ارزیابی مدل SCAD-LM-ANN با استفاده از پارامترهای آماری

یکی از روش‌های ارزیابی مدل توسعه یافته، محاسبه پارامترهای آماری مربوط به مدل می‌باشد. بنابراین پارامترهای آماری معرفی شده در بخش ۱-۵-۸-۴ برای فعالیت دارویی پیش‌بینی شده ترکیبات مجموعه آزمون و فعالیت‌های دارویی پیش‌بینی شده برای کل ترکیبات به روش رد مرحله‌ای تک تک محاسبه و در جدول ۶-۲ خلاصه شدند. نتایج حاصله در جدول ۶-۲ نشان می‌دهد که پارامترهای آماری در محدوده قابل قبول قرار دارند. بنابراین مطابق با بخش ۱-۵-۸-۴، بزرگ‌تر بودن مقادیر پارامترهای تروپشا و روی همچون  $R_0^2$ ،  $R_0^2$  نسبی،  $R_m^2$  و غیره از مقدار قابل قبول  $0/5$  و نزدیک بودن این پارامترها به مقدار  $R^2$  قدرت پیش‌بینی و تعمیم‌پذیری مدل توسعه یافته SCAD-LM-ANN اثبات می‌شود.

جدول ۲-۶ پارامترهای آماری محاسبه شده برای مجموعه آزمون و داده‌های پیش‌بینی شده با تکنیک LOO برای مدل SCAD-LM-ANN

ردیف	پارامتر آماری	LM-ANN		محدوده قابل قبول پارامتر آماری
		ترکیبات مجموعه آزمون	کل ترکیبات با تکنیک LOO	
		مقادیر pEC <sub>50</sub> پیش‌بینی شده با SCAD		
۱	PRESS	۱/۳۶	۶/۹	-
۲	SEP	۰/۳۵	۰/۳۵	-
۳	MAE	۰/۳۱	۰/۲۶	-
۴	REP(%)	۵/۷۹	۵/۶۳	-
۵	MSE	۰/۱۲	۰/۱۲	-
۶	MRE	۵/۱	۴/۳	-
۷	R <sup>2</sup>	۰/۹۲	-	R <sup>2</sup> > ۰/۶
۸	Q <sup>2</sup> <sub>LoO</sub>	-	۰/۸۱	Q <sup>2</sup> <sub>LoO</sub> > ۰/۵
۹	R <sup>2</sup> <sub>0</sub>	۰/۹۱	۰/۸۰	نزدیک به R <sup>2</sup>
۱۰	R <sup>2</sup> <sub>0</sub> نسبی	۰/۰۱	۰/۰۱	< ۰/۱
۱۱	R <sup>2</sup> <sub>m</sub>	۰/۹۲	۰/۸۱	> ۰/۵
۱۲	R <sup>2</sup> <sub>0</sub> '	۰/۸۹	۰/۸	نزدیک به R <sup>2</sup>
۱۳	R <sup>2</sup> <sub>0</sub> ' نسبی	۰/۰۳	۰/۰۱	< ۰/۱
۱۴	R <sup>2</sup> <sub>m</sub> '	۰/۹۲	۰/۸۱	> ۰/۵
۱۵	R-R	۰/۰۲	۰/۰۰	< ۰/۳
۱۶	k	۰/۹۶	۱/۰۰	۰/۸۵ ≤ k ≤ ۱/۱۵
۱۷	k'	۱/۰۴	۱/۰۰	۰/۸۵ ≤ k' ≤ ۱/۱۵

## ۲-۲-۴ ارزیابی مدل SCAD-LM-ANN با استفاده از دامنه کاربرد

دامنه کاربرد (AD) یک معیار مفید برای ارزیابی اعتمادپذیری<sup>۱</sup> مدل‌های QSAR است. دامنه

کاربرد مدل توسعه یافته SCAD-LM-ANN با استفاده از رسم نمودار ویلیام مورد تجزیه و تحلیل قرار

<sup>۱</sup>Reliability

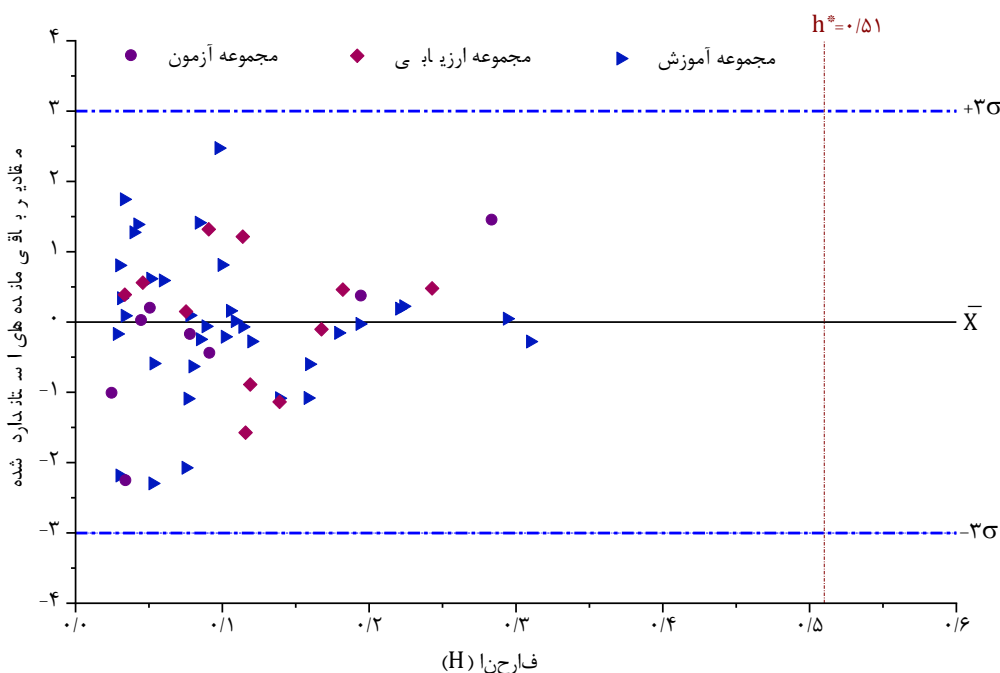
گرفت. طبق تعریف دامنه کاربرد در بخش ۱-۵-۸-۵، تعیین فضای ساختارهای شیمیایی برای بررسی پیش بینی قابل اعتماد مدل، ضروری است. به عبارت دیگر، درجه تعمیم پذیری یک مدل QSAR توسعه یافته به میزان گستردگی دامنه کاربرد وابسته است. در یک فضای شیمیایی با  $p$  توصیف کننده معین، تخمین هایی برای ساختارهای شیمیایی جدید با استفاده از داده های آموزشی به دست می آید [۱۴۲]. در این بخش به منظور بررسی دامنه کاربرد، محاسبه مقدار انحراف ( $H$ ) ضروری است. مقدار  $H$  یک ماده شیمیایی جدید، با اندازه گیری فاصله مالهونوبیس<sup>۱</sup> آن ترکیب از مرکز مجموعه آموزشی متناسب است. بنابراین ماتریس  $H$  با ابعاد  $۵۷ \times ۵۷$  با استفاده از رابطه ۱-۱۳ محاسبه شد. به این منظور ابتدا مقدار  $H$  مطابق با رابطه ۱-۱۳ محاسبه گردید و سپس بر اساس مقادیر  $pEC_{50}$  پیش بینی شده توسط مدل SCAD-LM-ANN، مقادیر باقی مانده های استاندارد شده محاسبه شد. از رسم نمودار باقی مانده های استاندارد شده بر حسب مقادیر  $H$ ، نمودار ویلیام استخراج شد و در شکل ۲-۷ نشان داده شد.

صحت آنالیز دامنه کاربرد، با قرارگیری داده های شیمیایی در دو محدوده اطمینان قابل قبول نمودار ویلیام مشخص می شود. اولین شرط قابل قبول، قرارگیری و عدم تجاوز مقدار باقی مانده های استاندارد شده داده های شیمیایی در محدوده ۳ برابری بزرگ تر / کوچک تری از انحراف استاندارد ( $\sigma$ ) می باشد. علاوه بر این شرط، مقادیر محاسبه شده  $H$  نیز نباید از مقدار حد آستانه یا حد هشدار  $h^*$  برابر با  $3p/n$  بزرگ تر باشند. در این معادله برابر با تعداد توصیف کننده های مدل به علاوه یک و  $n$  تعداد داده های مجموعه آموزش می باشد. بنابراین، قرارگیری داده های  $H$ ، در محدوده های قابل قبول، استحکام و قابل اعتماد بودن مدل های QSAR/QSPR توسعه یافته را اثبات می کند. با توجه به نتایج ارزیابی دامنه کاربرد نشان داده شده در نمودار ویلیام (شکل ۲-۷)، مشاهده می شود که مقادیر محاسبه شده برای  $H$  محاسبه شده کم تر از مقدار هشدار انحراف است و مقادیر باقی مانده های استاندارد محاسبه شده نیز در محدوده  $\pm 3\sigma$  می باشد. در نتیجه،

---

<sup>۱</sup>Mahalanobis

هیچ کدام از داده‌های به کار رفته در مدل SCAD-LM-ANN به‌عنوان داده دور افتاده شناخته نشده است و همه داده‌ها در دامنه کاربرد قابل قبول قرار گرفته‌اند.



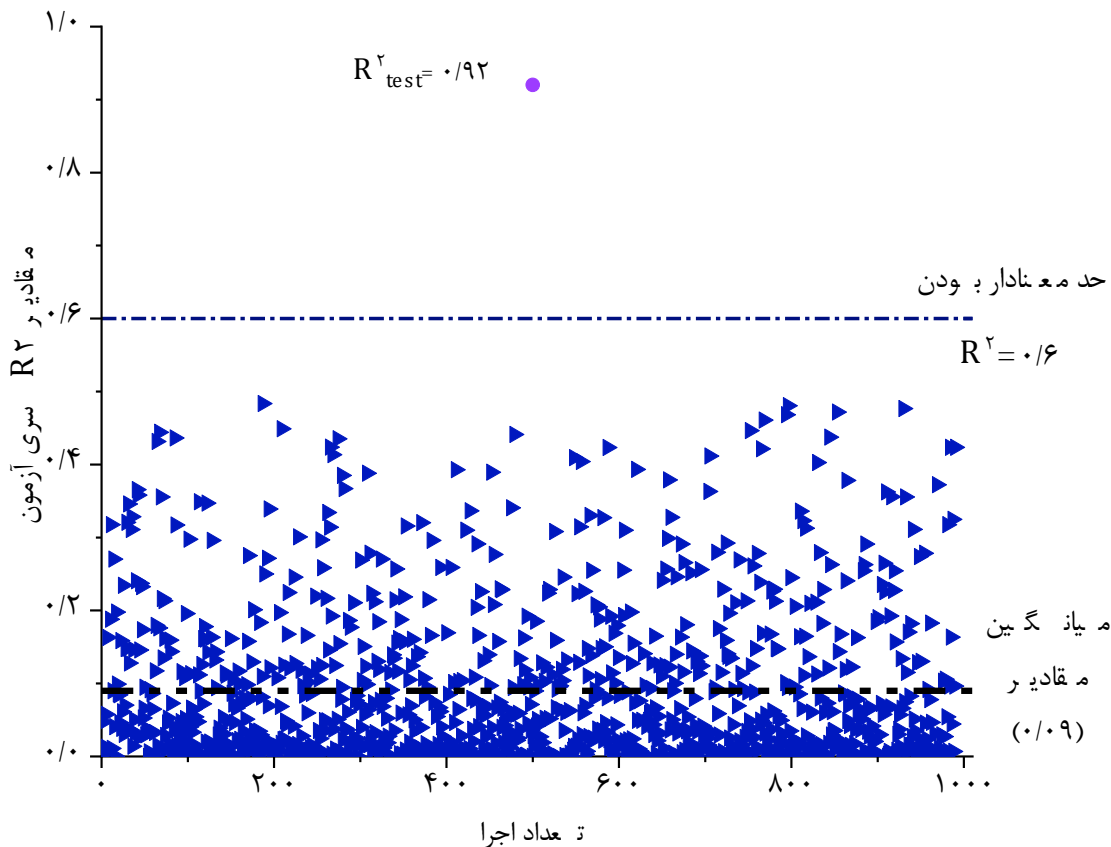
شکل ۲-۷ دامنه کاربرد مدل SCAD-LM-ANN، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $\pm 3\sigma$  و  $h^*$  است.

## ۲-۲-۵ ارزیابی مدل SCAD-LM-ANN با استفاده از آزمون Y-تصادفی

به منظور بررسی عدم وجود ارتباط تصادفی ایجاد شده توسط مدل SCAD-LM-ANN بین ویژگی‌های ساختار ترکیبات شیمیایی و فعالیت دارویی مربوطه از آزمون Y-تصادفی استفاده شد. برای انجام آزمون Y-تصادفی ابتدا مقادیر از این رو متغیر وابسته ( $pEC_{50}$ ) ترکیبات در مجموعه آموزش مورد مطالعه در محدوده حداقل و حداکثر مقادیر (۴/۶۲ تا ۷/۷۴) ۱۰۰۰ بار به صورت کاملاً تصادفی تغییر داده شدند و مدل شبکه عصبی با استفاده از داده‌های تصادفی متغیر وابسته توسعه داده شده و برای پیش بینی مقادیر  $pEC_{50}$  مجموعه آزمون به کار گرفته شدند. نتایج  $R^2$  حاصل از پیش بینی  $pEC_{50}$  ترکیبات مجموعه آزمون با ۱۰۰۰ مدل توسعه یافته با متغیر وابسته تصادفی،



در شکل ۸-۲ آورده شده است. میانگین مقادیر  $R^2$  برابر با  $0/09$  شد و مقادیر  $R^2$  حاصل برای هر بار اجرا از مقدار قابل قبول  $0/6$  کم تر است ( $R^2 < 0/6$ ). علاوه بر مقدار  $R^2$  مربوط به مجموعه آزمون که به وسیله مدل پیشنهادی SCAD-LM-ANN با استفاده از متغیر اصلی وابسته برابر با  $0/92$  می باشد که به طور معناداری با مقادیر  $R^2$  مجموعه آزمون پیش بینی شده در آزمون Y-تصادفی تفاوت دارد. می توان نتیجه گیری کرد که مدل پیشنهادی بر اساس ارتباط شانس بین متغیرهای مستقل و متغیر وابسته پایه گذاری نشده است و ارتباط بین ساختار شیمیایی (توصیف کننده های منتخب روش SCAD) و فعالیت ترکیبات مورد مطالعه در مدل توسعه یافته SCAD-LM-ANN به طور منطقی و معنادار است.



شکل ۸-۲ نمودار مقادیر  $R^2$  به دست آمده در آزمون Y-تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای Y-تصادفی و پیش بینی فعالیت ترکیبات مجموعه آزمون به وسیله مدل SCAD-LM-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه

## ۲-۳ پیش‌بینی فعالیت دارویی برخی از مشتقات ۳- chymotrypsin like

protease (3CL<sup>Pro</sup>) به‌عنوان بازدارنده‌های SARS-COV-2 با استفاده

### از مدل ALASSO-ANN

#### ۲-۳-۱ مقدمه

ویروس سندرم حاد تنفسی (SARS) می‌تواند منجر به سندرم حاد تنفسی شود و ویروس SARS Coronavirus (CoV) در سال ۲۰۰۳ در آسیا برای اولین بار شناسایی شد و بیش از ۸۰۰۰ نفر را دچار عفونت و ۸۰۰ نفر قربانی داشت [۱۴۳]. در پایان دسامبر ۲۰۱۹، SARS-CoV-2، به‌عنوان یک کرونا ویروس جدید، منجر به سندرم حاد تنفسی به‌نام کووید-۱۹ شد. همه‌گیری کووید-۱۹ به‌سرعت در سراسر جهان گسترش یافت و به یک خطر و چالش جهانی برای سلامت عمومی بشریت و ثبات اقتصادی همه جوامع تبدیل شد [۱۴۴، ۱۴۵]. سازمان جهانی بهداشت (WHO) این بیماری را در مارس ۲۰۲۰ به‌عنوان یک اپیدمی اعلام کرد و این بیماری تا به الان حدود ۲۵۸ میلیون عفونت و حدود ۵ میلیون قربانی داشته است. با توجه به همه‌گیری بیماری کووید-۱۹ و این واقعیت که تعداد بیماران و مرگ و میر روزانه در حال افزایش است، پیشنهاد ترکیبات جدید مشابه دارو برای پیشگیری یا درمان این بیماری یکی از مهم‌ترین چالش‌های محققان است. یک راه مناسب برای دستیابی به این ترکیبات، طراحی ترکیبات قوی جدید با مطالعه ساختار داروهایی است که بر بیماری‌های ناشی از ویروس‌هایی با ساختار مشابه ویروس کووید-۱۹ تأثیر می‌گذارند. در بین ویروس‌های خانواده SARS، ویروس‌های SARS-CoV دارای توالی ساختاری مشابه با SARS-CoV-2 هستند [۱۴۶]. کرونا ویروس جدید SARS-CoV-2 دارای ژنومی با حدود ۸۰ درصد شباهت به ویروس بتا درگیر در سندرم SARS-CoV است [۱۴۷]. مسدود کردن پروتئازهای ویروسی کدگذاری شده توسط ژنوم ویروس، نقش حیاتی در جلوگیری از تکثیر ویروس و پیشرفت بیماری ایفا می‌کند. در میان

پروتئازهای کرونا ویروس، پروتئاز papain-like ( $PL^{pro}$ ) و پروتئاز 3-chymotrypsin-like ( $3CL^{pro}$ ) برای تکثیر ویروس حیاتی هستند. آنزیم  $3CL^{pro}$  به عنوان یک هدف دارویی مناسب در مطالعات کووید-۱۹ توصیه می‌شود [۱۴۸]. در مطالعات اخیر SARS-CoV، مهارکننده‌های  $3CL^{pro}$  در دسته‌های متفاوتی، از جمله ایزاتین [۱۴۹, ۱۵۰]، پلی فنول‌ها [۱۵۱]، سینانسرین [۱۵۲]، پیریدین [۱۵۳]، سولفونیل دی بنزن [۱۵۴]، استامید [۱۵۵] و تانشینون [۱۵۶] و سایر مولکول‌های کوچک [۱۵۵, ۱۵۷, ۱۵۸] گنجانده شده‌اند. اگرچه مهارکننده‌های سنتز شده  $3CL^{pro}$  به دلیل گسترش بیش‌تر این بیماری در سطح جهان، مقاومت خوبی در برابر SARS-CoV نشان می‌دهند، اما همچنان برای مبارزه با عفونت این ویروس به ترکیباتی جدید با فعالیت ضد ویروسی نیاز است. یک راه سریع برای دستیابی به این ترکیبات، طراحی و کشف ترکیبات مشابه با مشتقات گزارش شده قبلی با فعالیت ضد ویروسی مناسب است. در این راستا، برخی از مشتقات مانند ایزاتین [۱۴۹, ۱۵۰]، پیریدین [۱۵۳]، سولفونیل دی بنزن [۱۵۴]، استامید [۱۵۵] و تانشینون [۱۵۶] با اثر بازدارندگی مناسب نسبت به SARS-CoV می‌توانند به عنوان کاندیدای مناسب در نظر گرفته شوند. روش‌های آزمون و خطای رایج برای طراحی دارو بسیار زمان‌بر و پرهزینه هستند. در نتیجه، رویکردهای طراحی دارویی به کمک کامپیوتر (CADD) برای طراحی ترکیبات جدید و قوی به شدت توصیه می‌شود. ابزارهای CADD و بیوانفورماتیک، به عنوان روش‌های سریع، ساده و اقتصادی [۱۵۹]، نقش مهمی را در طراحی ترکیبات بالقوه جدید دارند. یک مدل قدرتمند و معتبر QSAR دانش درستی را در مورد الگوهای ساختاری که همبستگی مناسبی با فعالیت دارویی ترکیبات شبه دارویی دارند، ارائه می‌کند [۴, ۱۶۲-۱۶۰]. مدل‌های QSAR دقت و صحت مطلوبی دارند و برای تخمین فعالیت دارویی ترکیبات جدید با ساختارهای مشابه با ساختار اصلی مجموعه داده‌ها مناسب هستند. پیش‌بینی فعالیت دارویی مهارکننده‌های جدید  $3CL^{pro}$  قبل از سنتز به محققان کمک می‌کند تا از سنتز ترکیبات با فعالیت دارویی کم و نامناسب اجتناب کنند و با استفاده از این رویکرد در طراحی دارو در زمان و هزینه صرفه‌جویی کنند. یکی از کاربردی‌ترین

جنبه‌های مطالعات QSAR که اخیراً مورد توجه بسیاری از محققان قرار گرفته است، استفاده از مدل‌های QSAR برای پیشنهاد ترکیبات جدید با فعالیت دارویی مناسب، بر اساس رابطه توسعه یافته بین فعالیت دارویی و توصیف‌کننده‌های مؤثر است. بنابراین این مفهوم مستلزم آن است که مدل QSAR با استفاده از تعداد کمی از توصیف‌کننده‌های مؤثر ساخته شده باشد. از این رو، استفاده از روش‌های انتخاب متغیر کارآمد و ترکیب آن‌ها با روش‌های مدل‌سازی قدرتمند برای پیشنهاد ترکیبات بالقوه ضروری است. روش‌های انتخاب متغیر با کارایی بالا، توصیف‌کننده‌های اضافی را حذف می‌کنند، درحالی‌که توصیف‌کننده‌های مؤثر مرتبط با فعالیت دارویی را تا حد امکان حفظ می‌کنند. در مطالعات QSAR، پرداختن به داده‌های با ابعاد بالا و حذف توصیف‌کننده‌های اضافی و بی‌اهمیت همیشه موضوع چالش برانگیز محققان بوده است [۱۶۲]. از بین روش‌های انتخاب متغیر، روش‌های جریمه‌ای به دلیل داشتن مزایایی چون، بایاس کم، تنگی مناسب، پایداری و عملکرد بالا در حضور هم‌خطی بسیار مورد توجه قرار گرفته‌اند. از جمله روش‌های رگرسیون جریمه‌شده، می‌توان به لاسو تطبیقی (ALASSO) اشاره کرد که برای غلبه بر محدودیت‌های ذاتی روش‌های کلاسیک مناسب است [۳۰]. بنابراین، در این بخش از مطالعه، برای ساخت مدل QSAR از روش انتخاب متغیر ALASSO برای انتخاب مؤثرترین توصیف‌کننده‌ها استفاده شد، که در ادامه به آن پرداخته می‌شود. در این مطالعه، یک مدل QSAR تنک با قدرت پیش‌بینی و تفسیرپذیری بالا برای پیش‌بینی فعالیت دارویی ترکیبات ضد CoV-2 ارائه شد. به این منظور، ترکیبی از ALASSO و روش مدل‌سازی شبکه عصبی برای مدل‌سازی غیرخطی فعالیت دارویی برخی از بازدارنده‌های 3CL<sup>pro</sup> ساخته شد. بنابراین، مدل QSAR پیشنهادی در این مطالعه، مهارکننده‌های جدید 3CL<sup>pro</sup> با خواص فارماکوکینتیک مناسب را پیشنهاد می‌کند.

## ۲-۳-۲ مجموعه داده‌ها

مجموعه داده‌ها شامل ۹۰ مهارکننده 3CL<sup>pro</sup> می‌باشد که ساختار ترکیبات و فعالیت دارویی (IC<sub>50</sub>) آن‌ها از مقالات منتشر شده استخراج شد [۱۴۹, ۱۵۰, ۱۵۳-۱۵۶]. لگاریتم مقادیر (1/IC<sub>50</sub>) محاسبه شد و به‌عنوان پاسخ در مراحل مدل‌سازی QSAR مورد استفاده قرار گرفت. IC<sub>50</sub> غلظت مؤثر (بر حسب مولار) مورد نیاز برای بازداری ساختارهای شیمیایی ۹۰ ترکیب به یک فایل فرمت سیستم ورودی خطی ورودی مولکولی ساده شده (SMILES) تبدیل شد و همراه با مقادیر pIC<sub>50</sub> در جدول ۲-۷ خلاصه شد. تقسیم‌بندی مجموعه داده‌ها با استفاده از الگوریتم KS انجام شد و داده‌ها به سه دسته آموزش (۶۴ ترکیب)، ارزیابی (۱۳ ترکیب) و آزمون (۱۳ ترکیب) تقسیم شدند و در مراحل مختلف مدل‌سازی QSAR مورد استفاده قرار گرفتند.

---

<sup>1</sup>Simplified molecular-input line-entry system

جدول ۷-۲ مجموعه داده‌ها به همراه مقادیر واقعی و پیش‌بینی شده pIC<sub>50</sub>

ردیف	SMILES	pIC <sub>50</sub> واقعی	pIC <sub>50</sub> پیش‌بینی شده
۱ <sup>v</sup>	<chem>c1(ccc2c(c1)C(=O)C(=O)N2Cc1c(oc1C)C)C#N</chem>	۵/۱۴	۵/۵۸
۲	<chem>c1(ccc2c(c1)C(=O)C(=O)N2Cc1c(cc(cc1)F)Cl)I</chem>	۵/۰۳	۵/۳
۳ <sup>t</sup>	<chem>c1(ccc2c(c1)C(=O)C(=O)N2C[C@@H]1COc2c(O1)cccc2)I</chem>	۴/۸۷	۴/۷۹
۴	<chem>c1ccc2c(c1)C(=O)C(=O)N2Cc1cc2c(s1)cccc2</chem>	۴/۸۸	۵/۱۲
۵ <sup>v</sup>	<chem>c1cc(c2c(c1)C(=O)C(=O)N2Cc1cc2c(s1)cccc2)N(=O)=O</chem>	۵/۷	۵/۵۸
۶ <sup>t</sup>	<chem>c1cc(c2c(c1)C(=O)C(=O)N2Cc1cc2c(s1)cccc2)Br</chem>	۶/۰۱	۶/۷۸
۷	<chem>c1(ccc2c(c1)C(=O)C(=O)N2Cc1cc2c(s1)cccc2)F</chem>	۵/۳۲	۵/۱۹
۸	<chem>c1(ccc2c(c1)C(=O)C(=O)N2Cc1cc2c(s1)cccc2)I</chem>	۶/۰۲	۵/۱۹
۹	<chem>c1ccc2c(c1Cl)C(=O)C(=O)N2Cc1cc2c(s1)cccc2</chem>	۴/۹۵	۵/۶۳
۱۰	<chem>c1(ccc2c(c1)C(=O)C(=O)N2C/C=C/c1cc2c(s1)cccc2)I</chem>	۴/۶۳	۴/۶۴
۱۱	<chem>c1(ccc2c(c1)C(=O)C(=O)N2Cc1ccc(s1)C(=O)Nc1ccc(cc1)Cl)I</chem>	۴/۹	۴/۹۷
۱۲	<chem>c1(ccc2c(c1)C(=O)C(=O)N2Cc1ccc(s1)C(=O)N1CCCCC1)I</chem>	۴/۷۶	۴/۸۱
۱۳	<chem>c1c(c(cc(c1Cl)S(=O)(=O)c1c(cc(cc1N(=O)=O)C(F)(F)F)N(=O)=O)C)Cl</chem>	۶/۵۲	۶/۲۶
۱۴	<chem>c1c(ccc(c1)S(=[O-])(=[O-])c1ccc(cc1)OC(=O)C(=C(Cl)Cl)Cl)OC(=O)C(=C(Cl)Cl)Cl</chem>	۶/۰۵	۵/۹۶
۱۵	<chem>c1c(ccc(c1)[S+3])(=[O-])(=[O-])[C-]1=CN=C(N=C1N)N)Cl</chem>	۵/۲۲	۵/۳۴
۱۶	<chem>c1c(ccc(c1)[S@@])(=[O-])(=O)c1ccc(cc1N(=O)=O)C(F)(F)F)Cl</chem>	۴/۹۲	۴/۷۶
۱۷	<chem>c1cccc(c1)[S@@])(=[O-])(=O)c1nc(c(c1C#N)C)N(=O)=O)C</chem>	۴/۸۹	۵/۱۷
۱۸	<chem>c1c(ccc(c1)S(=[O-])(=[O-])c1nc(c(c1C#N)C)N(=O)=O)C)C</chem>	۴/۸۹	۴/۸۱
۱۹ <sup>t</sup>	<chem>c1c(ccc(c1)[S@@])(=O)(=[O-])c1nc(cc1)N(=O)=O)Cl</chem>	۴/۸۲	۵/۰۶
۲۰	<chem>c1(ccc(cc1)N/C=C(/C(=O)OCC)\C#N)S(=[O-])(=[O-])c1ccc(cc1)N/C=C(/C(=O)OCC)\C#N</chem>	۴/۸	۴/۷۷
۲۱ <sup>v</sup>	<chem>c1(cc(c(c(c1)Br)O)Br)S(=[O-])(=[O-])c1ccc(cc1)C(=O)O</chem>	۴/۸	۴/۹۶
۲۲	<chem>c1(ccccc1C(=O)C)S(=[O-])(=[O-])c1c(cccc1)C(=O)O</chem>	۴/۸	۴/۵۸
۲۳ <sup>t</sup>	<chem>c1(ccc(cc1)N(=O)=O)S(=[O-])(=[O-])c1ccc(cc1)N(=O)=O</chem>	۴/۶	۴/۴۴
۲۴	<chem>c1(ccc(cc1)NC=C(C(=O)OCC)C(=O)OCC)S(=[O-])(=[O-])c1ccc(cc1)NC=C(C(=O)OCC)C(=O)OCC</chem>	۴/۴۹	۴/۶۵
۲۵ <sup>v</sup>	<chem>c1cccc(c1)C#Cc1oc(cc1)C(=O)Sc1nnc([nH]1)C(F)F)F</chem>	۵/۵۲	۵/۰۳
۲۶	<chem>c1(cc(n1C)C(F)(F)F)c1ccc(s1)C(=O)Nc1c(cc1)C)N(=O)=O</chem>	۵/۳	۴/۸۱
۲۷ <sup>v</sup>	<chem>c1ccsc1[S+3])(=[O-])(=[O-])Nc1[nH]nc(c1)c1cc(oc1C)C(C)C)C</chem>	۵/۰۰	۴/۸
۲۸	<chem>c1ccsc1C(=O)Nc1[nH]nc([c-]1[S+3])(=[O-])(=[O-])c1ccc(cc1)Cl)SC</chem>	۴/۸۲	۴/۷۶
۲۹ <sup>t</sup>	<chem>c1coc(c1)CNc1sc(c2c1C(=O)CC(C2)(C)C)C#N</chem>	۴/۸	۴/۶۷
۳۰	<chem>c1ccc(s1)c1n(nc(c1C#N)SC)c1c(c(n1C)C)N(=O)=O</chem>	۴/۷۴	۴/۷۲
۳۱	<chem>c1esc(c1)c1ccc(s1)CNc1n(nc(c1N(=O)=O)C)C</chem>	۴/۷	۴/۶۶
۳۲	<chem>c1esc(n1)C(=O)CC1=NCCS1</chem>	۴/۴	۴/۶۹
۳۳	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCN(CC1)C</chem>	۴/۱۱	۴/۳۹
۳۴ <sup>v</sup>	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCN(CC1)Cc1cc(ccc1)Cl</chem>	۴/۵۰	۴/۴۳
۳۵	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCN(CC1)Cc1cc(c(c1)OC)OC)OC</chem>	۴/۴۹	۴/۴
۳۶	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCN(CC1)CCc1cccc1</chem>	۴/۴۶	۴/۳۸
۳۷	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCN(CC1)C(=O)c1ccc1</chem>	۵/۰۰	۴/۹۵
۳۸	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCN(CC1)c1ccc1</chem>	۴/۲۹	۴/۷۱

ادامه جدول ۲-۷

ردیف	SMILES	pIC <sub>۵۰</sub> واقعی	pIC <sub>۵۰</sub> پیش‌بینی شده
۳۹ <sup>v</sup>	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCCCC1</chem>	۵/۳۵	۵
۴۰	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCOCC1</chem>	۴/۹۰	۴/۸۳
۴۱	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCC(CC1)C</chem>	۵/۹۳	۵/۵۹
۴۲ <sup>v</sup>	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1CCCC[C@@H]1C</chem>	۵/۶۵	۴/۹۱
۴۳	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)S(=[O-])(=[O-])N1C[C@H](C[C@H](C1)C)C</chem>	۵/۳۷	۵/۵۲
۴۴	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)C)S(=[O-])(=[O-])N1CCN(CC1)C</chem>	۴/۹۳	۴/۳۹
۴۵	<chem>c1(ccc2c(c1)C(=O)C(=O)N2)Cc1cccc1)S(=[O-])(=[O-])N1CCN(CC1)C</chem>	۴/۱۷	۴/۴۵
۴۶	<chem>c12c(N(C(=O)C1=O)C)Cc1ccc3c(c1)cccc3)ccc(c2)S(=[O-])(=[O-])N1CCN(CC1)C</chem>	۴/۰۸	۴/۴۱
۴۷	<chem>c12c(N(C(=O)C1=O)C)ccc(c2)S(=[O-])(=[O-])N1CCN(CC1)CCc1cccc1</chem>	۴/۸۶	۴/۳۸
۴۸	<chem>c12c(N(C(=O)C1=O)C)ccc(c2)S(=[O-])(=[O-])N1CCN(CC1)c1cccc1</chem>	۵/۲۶	۴/۷۱
۴۹ <sup>v</sup>	<chem>c12c(N(C(=O)C1=O)C)Cc1cccc1)ccc(c2)S(=[O-])(=[O-])N1CCCCC1</chem>	۴/۸۵	۴/۹۸
۵۰	<chem>c12c(N(C(=O)C1=O)C)ccc(c2)S(=[O-])(=[O-])N1CCOCC1</chem>	۵/۰۰	۵/۰۵
۵۱ <sup>t</sup>	<chem>c12c(N(C(=O)C1=O)C)Cc1cccc1)ccc(c2)S(=[O-])(=[O-])N1CCOCC1</chem>	۴/۸۶	۴/۷۵
۵۲	<chem>c12c(N(C(=O)C1=O)C)Cc1ccc3c(c1)cccc3)ccc(c2)S(=[O-])(=[O-])N1CCOCC1</chem>	۴/۴۰	۴/۸۶
۵۳	<chem>c12c(N(C(=O)C1=O)C)ccc(c2)S(=[O-])(=[O-])N1CCC(CC1)C</chem>	۵/۹۸	۵/۵۷
۵۴ <sup>t</sup>	<chem>c12c(N(C(=O)C1=O)C)Cc1cccc1)ccc(c2)S(=[O-])(=[O-])N1CCC(CC1)C</chem>	۵/۷۷	۴/۹۳
۵۵	<chem>c12c(N(C(=O)C1=O)C)Cc1ccc3c(c1)cccc3)ccc(c2)S(=[O-])(=[O-])N1CCC(CC1)C</chem>	۴/۷۵	۵/۰۵
۵۶	<chem>c12c(N(C(=O)C1=O)C)ccc(c2)S(=[O-])(=[O-])N1C[C@H](C[C@H](C1)C)C</chem>	۵/۵۵	۵/۵۵
۵۷	<chem>c12c(N(C(=O)C1=O)C)Cc1cccc1)ccc(c2)S(=[O-])(=[O-])N1C[C@H](C[C@H](C1)C)C</chem>	۵/۳۳	۵/۳۹
۵۸	<chem>c1c(cc(cc1Cl)NC(=O)C)Sc1nccc(n1)c1nc(sc1)c1cccc1)Cl</chem>	۵/۵۲	۵/۶۷
۵۹	<chem>c1c(ccc(c1)NC(=O)C)Sc1nc(ccn1)c1ccc(s1)c1cc(n(n1)C)C(F)(F)F)Cl</chem>	۵/۰۰	۴/۶۷
۶۰ <sup>v</sup>	<chem>c1cc(ccc1NC(=O)C)Sc1nc(ccn1)c1c(c(sc1SC)c1nc(es1)C)C)Cl</chem>	۴/۹۶	۴/۶۹
۶۱ <sup>t</sup>	<chem>c1(ccccc1NC(=O)C)Sc1nc(ccn1)c1c(c(sc1SC)c1nc(es1)C)C)Cl</chem>	۴/۹۲	۴/۸۶
۶۲	<chem>c1c(ccc(c1)NC(=O)C)Sc1nccc(n1)c1nc(sc1)c1nc(sc1)C)Cl</chem>	۴/۸۵	۵/۰۱
۶۳ <sup>t</sup>	<chem>c1c(ccc(c1)NC(=O)C)Sc1nccc(n1)c1cc(no1)c1cccc1)Cl</chem>	۴/۸۲	۴/۵
۶۴	<chem>c1cc(ccc1NC(=O)C)Sc1nc(ccn1)c1cc(no1)c1ccc(cc1Cl)Cl)C(F)(F)F</chem>	۴/۸۲	۴/۸۱
۶۵ <sup>t</sup>	<chem>c1cc(ccc1NC(=O)C)Sc1nccc(n1)c1cc(no1)c1c(cccc1)Cl)Cl</chem>	۴/۸۲	۴/۵۱
۶۶	<chem>c1cc(ccc1NC(=O)C)Sc1nc(cc(n1)O)CCC)Cl</chem>	۴/۵۲	۴/۴۹
۶۷	<chem>c1c(ccc(c1)NC(=O)C)Sc1nc(c(c(=O)[nH]1)C#N)c1cccc(c1)OC)S(=[O-])(=[O-])[NH2-]</chem>	۴/۴۰	۴/۵۳
۶۸	<chem>c1c(ccc(c1)NC(=O)C)Sc1nc(ccn1)c1cccs1)C(C)C</chem>	۴/۴۰	۴/۴۹
۶۹	<chem>n1c(nc(cc1c1cccc1)O)SCC(=O)Nc1ccc(c(c1)OC)OC</chem>	۴/۳۵	۴/۶۱
۷۰ <sup>t</sup>	<chem>c1cc(ccc1NC(=O)C)Sc1[nH]c(=O)c(c(n1)c1cc(ccc1)OC)C#N)C(=O)C</chem>	۴/۲۲	۴/۴۱

ادامه جدول ۷-۲

ردیف	SMILES	pIC <sub>50</sub> . واقعی	pIC <sub>50</sub> . پیش‌بینی شده
۷۱	<chem>c1(cccc(c1)NC(=O)[C@@H](CC)Sc1nc(c(c(=O)[nH]1)C#N)c1c ccc(c1)OC)C(=O)C</chem>	۴/۲۲	۴/۴۱
۷۲	<chem>c1(ccc(c(c1)NC(=O)CSc1[nH]c(=O)cc(n1)C)Oc1cccc1)Cl</chem>	۴/۰۰	۴/۴۶
۷۳	<chem>c1c(ccc(c1)c1oc(cc1)C(=O)Oc1ncc(c1)Cl)Cl</chem>	۷/۲۰	۶/۵۲
۷۴ <sup>t</sup>	<chem>c1c(ccc(c1)c1oc(cc1)C(=O)Oc1ncc(c1)Cl)N(=O)=O</chem>	۷/۲۲	۶/۶۷
۷۵ <sup>v</sup>	<chem>c1c(cc(c(c1)c1oc(cc1)C(=O)Oc1ncc(c1)Cl)N(=O)=O)Cl</chem>	۶/۹۱	۶/۷۴
۷۶	<chem>c1ccc(c(c1)c1oc(cc1)C(=O)Oc1ncc(c1)Cl)N(=O)=O</chem>	۶/۶۸	۶/۷۸
۷۷	<chem>c1cc(cc(c1)c1oc(cc1)C(=O)Oc1ncc(c1)Cl)N(=O)=O</chem>	۶/۳۰	۶/۷۷
۷۸	<chem>c1(C(=O)Oc2ncc(c2)Cl)ccncc1</chem>	۶/۷۹	۶/۶۳
۷۹	<chem>c1(C(=O)Oc2ncc(c2)Cl)enccc1</chem>	۶/۱۶	۶/۳۹
۸۰	<chem>c1(C(=O)Oc2ncc(c2)Cl)ccc(cc1)Cl</chem>	۶/۳۶	۶/۴۱
۸۱ <sup>v</sup>	<chem>c1(C(=O)Oc2ncc(c2)Cl)c(cccc1)N(=O)=O</chem>	۶/۴۸	۶/۱۳
۸۲	<chem>c1(C(=O)Oc2ncc(c2)Cl)cc(ccc1)N(=O)=O</chem>	۶/۱۶	۶/۰۱
۸۳ <sup>t</sup>	<chem>c1(C(=O)Oc2ncc(c2)Cl)c2c(ccc1)cccc2</chem>	۶/۹۱	۶/۶۶
۸۴ <sup>v</sup>	<chem>c1(C(=O)Oc2ncc(c2)Cl)c(=O)oc2c(c1)cccc2</chem>	۶/۹۷	۶/۷۸
۸۵	<chem>C1CC(c2c(C1)c1c(cc2)c2c(C(=O)C1=O)c(co2)C)(C)C</chem>	۴/۰۵	۴/۳۹
۸۶	<chem>C1C[C@](c2c(C1)c1c(cc2)c2c(C(=O)C1=O)c(co2)C)(C)CO</chem>	۴/۶۱	۴/۴۲
۸۷	<chem>C1C[C@](c2c(C1)c1c(cc2)c2c(C(=O)C1=O)c(co2)C)(C)C(=O)O C</chem>	۴/۶۸	۴/۴۲
۸۸	<chem>c1cc(c2c(c1)c1c(cc2)c2c(C(=O)C1=O)c(co2)C)C</chem>	۴/۴۱	۴/۵۸
۸۹	<chem>c1cc(c2c(c1)c1c(cc2)C2=C(C(=O)C1=O)[C@H](CO2)C)C</chem>	۴/۸۴	۴/۸۷
۹۰	<chem>C1CC(c2c(C1)c1c(cc2)C=C(C(=O)C1=O)C(C)C)(C)C</chem>	۴/۶۸	۴/۴۴



## ۲-۳-۳ رسم و بهینه‌سازی ساختار بازدارنده‌های 3CL<sup>pro</sup>

ساختارهای سه بعدی ۹۰ بازدارنده 3CL<sup>pro</sup> با استفاده از نرم‌افزار هایپرکم و مطابق با بخش ۱-۵-۳ رسم شد. ساختارهای شیمیایی مطابق با روش کار بخش ۱-۵-۳ تا رسیدن به حداقل انرژی بهینه شدند و در نهایت با فرمت \*.hin ذخیره شدند.

## ۲-۳-۴ استخراج توصیف‌کننده‌ها

ساختارهای بهینه ۹۰ بازدارنده 3CL<sup>pro</sup> مورد مطالعه در نرم‌افزار دراگون فراخوانی شدند و سپس برای هر ترکیب شیمیایی به تعداد ۳۲۲۴ توصیف‌کننده محاسبه شدند.

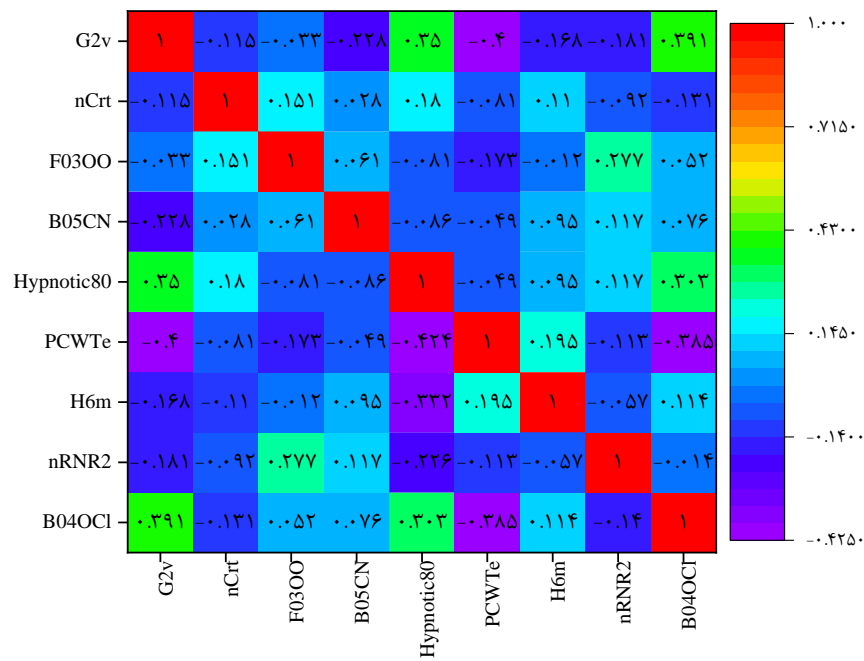
## ۲-۳-۵ پیش‌پردازش و انتخاب توصیف‌کننده‌های مؤثر

توصیف‌کننده‌های مولکولی با استفاده از ساختارهای بهینه و با استفاده از نرم‌افزار دراگون محاسبه شدند. از بین توصیف‌کننده‌های محاسبه شده، توصیف‌کننده‌هایی با مقادیر ثابت و نسبتاً ثابت (توصیف‌کننده‌هایی با واریانس کم‌تر از ۰/۰۰۱) با اجرای بسته نرم‌افزاری caret در برنامه R حذف شدند [۱۶۳]. بنابراین حدود ۱۳۸۶ توصیف‌کننده که اطلاعات خاصی به مدل اضافه نمی‌کنند، از مجموعه داده‌ها حذف شدند. از بین دو توصیف‌کننده با همبستگی بالاتر از ۰/۹، توصیف‌کننده‌ای که همبستگی کم‌تری با پاسخ داشت حذف شد. پس از انجام مراحل پیش‌پردازش انجام شده، ماتریسی شامل ۶۴۵ توصیف‌کننده برای مطالعات بیش‌تر ذخیره شد. سپس روش انتخاب متغیر ALASSO به‌عنوان روش انتخاب متغیر جریمه‌شده بر روی ماتریس داده‌ها (شامل ۶۴۵ متغیر مستقل و  $pIC_{50}$  به‌عنوان متغیر وابسته) اجرا شد. رگرسیون ALASSO با روش ارزیابی تقاطعی ۱۰ فولد موجود در بسته نرم‌افزاری parcor در برنامه R بر روی مجموعه داده‌های ارزیابی و آموزش اجرا شد تا مؤثرترین توصیف‌کننده‌ها انتخاب شود [۳۴]. با اجرای روش ALASSO به تعداد ۹ توصیف‌کننده مربوط به  $\lambda_{min}$  ( $\lambda$  دارای کم‌ترین خطای ارزیابی تقاطعی) به‌عنوان مؤثرترین توصیف‌کننده‌ها انتخاب شدند. مقادیر مربوط به ضرایب رگرسیونی ۹ توصیف‌کننده منتخب روش

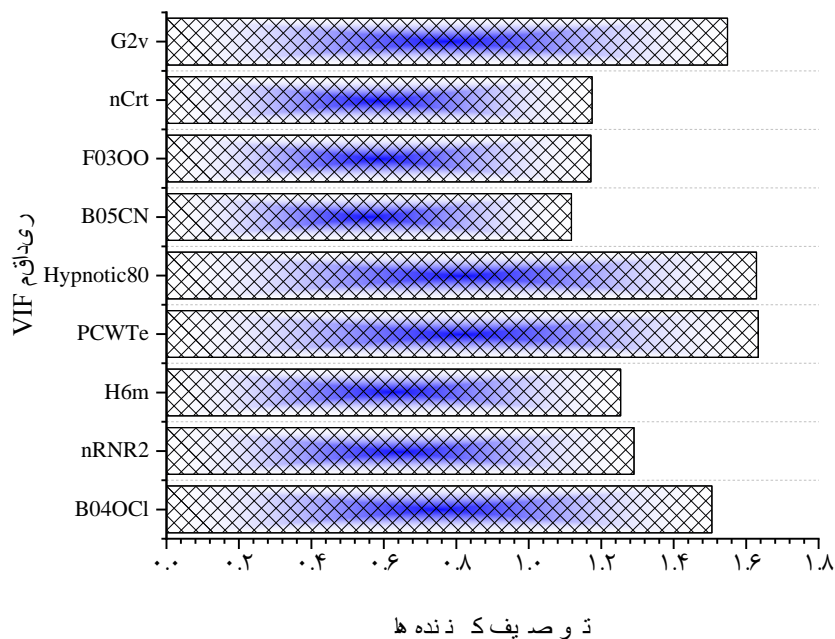
ALASSO در جدول ۸-۲ آورده شده است. ارزیابی‌های مربوط به کیفیت آماری توصیف‌کننده‌های منتخب روش ALASSO، از جمله عدم وجود همبستگی و هم‌خطی بین توصیف‌کننده‌ها با محاسبه مقادیر ضرایب همبستگی بین دو توصیف‌کننده و مقادیر افزایش تورم واریانس (VIF) توصیف‌کننده (مطابق با رابطه ۱۰-۱) مطالعه شد. نتایج مربوط به آن‌ها در نمودارهای نقشه رنگی و VIF (شکل ۹-۲ و شکل ۱۰-۲) آورده شده است. شکل ۹-۲ نشان می‌دهد که همبستگی معناداری بین توصیف‌کننده‌های انتخاب شده با روش ALASSO وجود ندارد. علاوه بر این، در شکل ۱۰-۲ نیز، مقادیر کم‌تر از ۱۰ مربوط به آنالیز افزایش تورم واریانس مشاهده می‌شود. نتایج حاکی از آن است که بین توصیف‌کننده‌های منتخب روش ALASSO هم‌خطی شدیدی وجود ندارد [۱۳۸، ۱۳۹].

جدول ۸-۲ توصیف‌کننده‌های منتخب ALASSO

ردیف	نماد	طبقه‌بندی	معنا	ضرایب استاندارد شده ALASSO
۱	B04[O-Cl]	2D binary fingerprints	Presence/absence of O - Cl at topological distance 4	۰/۶۶
۲	nRNR2	Functional group counts	number of tertiary amines (aliphatic)	۰/۱۲
۳	H6m	GETAWAY descriptors	H autocorrelation of lag 6 / weighted by mass	۰/۱
۴	PCWTe	Charge descriptors	partial charge weighted topological electronic index	۰/۰۹
۵	Hypnotic-80	Drug-like indices	Ghose-Viswanadhan-Wendoloski hypnotic-like index at 80%	۰/۰۸
۶	B05[C-N]	2D binary fingerprints	Presence/absence of C - N at topological distance 5	۰/۰۷
۷	F03[O-O]	2D frequency fingerprints	Frequency of O - O at topological distance 3	۰/۰۴
۸	nCrt	Functional group counts	number of ring tertiary C(sp3)	۰/۰۳
۹	G2v	WHIM descriptors	nd component symmetry directional ۲ WHIM index / weighted by van der Waals volume	۰/۰۰۲



شکل ۹-۲ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب روش ALASSO



شکل ۱۰-۲ نمودار مقادیر VIF توصیف‌کننده‌های منتخب روش ALASSO

## ۲-۳-۶ مدل سازی شبکه عصبی با استفاده از توصیف کننده های منتخب ALASSO

رابطه بین توصیف کننده های ساختاری و  $pIC_{50}$  با استفاده از یک مدل شبکه عصبی مصنوعی پیشخور با الگوریتم آموزشی پس انتشار خطا ایجاد شد. همان طور که در بخش ۲-۳-۶ اشاره شد، مجموعه داده ها به سه دسته آموزش، ارزیابی و آزمون تقسیم شد. گروه اول (۶۴ ترکیب) برای آموزش مدل شبکه عصبی و از گروه دوم (۱۳ ترکیب) برای ارزیابی و انتخاب مدل بهینه استفاده شد. باقی مانده ترکیبات در مدل سازی شبکه عصبی شرکت نداشتند و به عنوان یک مجموعه آزمون خارجی مستقل، برای بررسی تعمیم پذیری مدل در شبکه عصبی تعریف شدند. به منظور یافتن شبکه عصبی بهینه، پارامترهای شبکه عصبی از جمله تعداد ورودی ها، تعداد گره های لایه پنهان، دوره های آموزشی و توابع انتقال (لگاریتم سیگموئیدی و تانژانت هایپربولیک سیگموئیدی به ترتیب با توابع  $\logsig$  و  $tansig$  در جعبه ابزار متلب) و توابع آموزش لونبرگ مارکورات و تنظیم بایزین (به ترتیب با توابع آموزشی  $trainlm$  و  $trainbr$  در جعبه ابزار متلب شناخته می شوند) به طور هم زمان بهینه شد. در تمام مدل های شبکه عصبی مصنوعی، تابع خطی به عنوان لایه خروجی استفاده شد. توصیف کننده های منتخب روش ALASSO بر اساس بزرگی ضرایب استاندارد شده ALASSO (جدول ۲-۸) چیده شدند و به عنوان ورودی در ساخت مدل شبکه عصبی مورد استفاده قرار گرفتند. در بهینه سازی هم زمان پارامترهای شبکه عصبی مصنوعی، تعداد توصیف کننده ها، گره های لایه پنهان و دوره های آموزشی به طور هم زمان از ۲ تا ۹ (با گام ۱)، ۲ تا ۱۰ (با گام ۱) و ۵ تا ۵۰ (با گام ۵) تغییر یافت. مدل های ممکن ایجاد شده ANN، با استفاده از مجموعه داده های آموزش (شامل ۶۴ ترکیب) ساخته شدند و برای پیش بینی مقادیر  $pIC_{50}$  مربوط به ترکیبات مجموعه ارزیابی استفاده شدند. مقادیر MSE مجموعه ارزیابی برای تمام مدل های آموزش دیده محاسبه شد و به عنوان معیاری برای انتخاب بهترین مدل ANN مورد استفاده قرار گرفت. نتایج معماری های شبکه و مقادیر MSE برای چهار مدل ANN با توابع آموزش و انتقال متفاوت که دارای کم ترین مقدار MSE هستند، در جدول ۲-۹ آورده

شده است. نتایج نشان می‌دهند که مدل ANN با ۹ توصیف کننده منتخب روش ALASSO به‌عنوان ورودی، ۲ گره در لایه پنهان و ۵ دور آموزش با تابع انتقال تانژانت هایپربولیک سیگموئیدی و الگوریتم آموزشی LM دارای کم‌ترین مقدار MSE برای پیش بینی داده‌های مجموعه ارزیابی می‌باشد. بنابراین مدل با نماد ALASSO-LM-ANN نشان داده شده با معماری ۱-۲-۹ به‌عنوان مدل برتر برای پیش‌بینی مقادیر  $pIC_{50}$  ترکیبات مورد مطالعه انتخاب شد.

جدول ۲-۹ ساختارهای شبکه‌های توسعه یافته با توصیف‌کننده‌های منتخب ALASSO با کمترین MSE مجموعه ارزیابی

تعداد توصیف کننده	تابع آموزش	تابع انتقال	تعداد گره	تعداد دور آموزش	MSE <sub>validation</sub>	R <sup>2</sup> <sub>validation</sub>
۹	تنظیم بایزین	لگاریتم-سیگموئید	۵	۵	۰/۱۴	۰/۸۲
۹	لونبرگ-مارکوارت	لگاریتم-سیگموئید	۲	۴۰	۰/۱۲	۰/۸۵
۹	تنظیم بایزین	تانژانت-سیگموئید	۹	۵	۰/۱۶	۰/۸۱
۹	لونبرگ-مارکوارت	تانژانت-سیگموئید	۲	۵	۰/۱۱	۰/۸۶

با هدف مقایسه عملکرد ALASSO در انتخاب متغیر برای مدل‌سازی ANN، از روش کلاسیک متداول SR برای انتخاب موثرترین توصیف کننده‌ها استفاده شد. پس از اجرای SR روی مجموعه داده‌های ارزیابی، ۱۸ توصیف کننده به‌عنوان بهترین زیرمجموعه از توصیف کننده‌های SR انتخاب شد. به منظور مقایسه کارایی SR و SCAD، به تعداد ۹ توصیف کننده روش SR (چیده شده بر اساس بزرگی ضرایب رگرسیونی) برای طراحی و بهینه‌سازی مدل ANN مورد استفاده قرار گرفتند. نتایج نشان داد مدل ANN با تابع آموزش LM و با استفاده از زیر مجموعه‌های ۲ تا ۹ تایی از توصیف کننده‌های منتخب روش SR با تعداد ۳ گره در لایه پنهان و ۵ دور آموزش حداقل MSE را برای مجموعه ارزیابی ایجاد نمود. بنابراین مدل بهینه SR-LM-ANN برای پیش بینی فعالیت دارویی بازدارنده‌های 3CL<sup>pro</sup> موجود در مجموعه آزمون استفاده شد.

## ۲-۳-۷ ارزیابی مدل ALASSO-LM-ANN

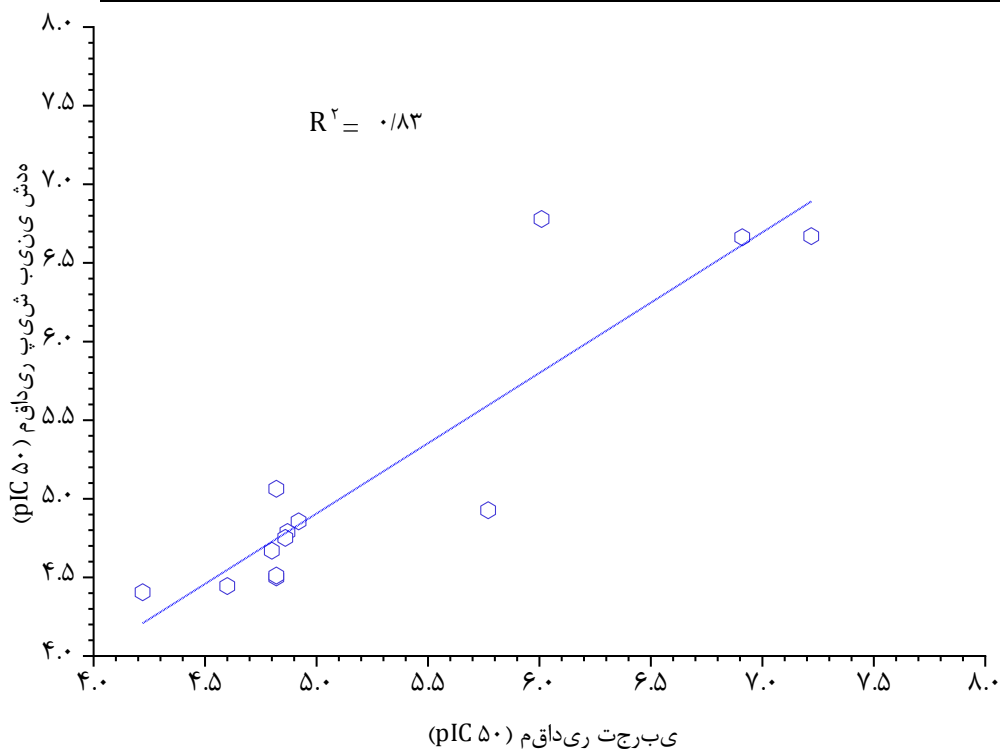
ارزیابی قدرت پیش‌بینی مدل شبکه عصبی پیشنهادی (ALASSO-LM-ANN) گامی مهم در مطالعات QSAR است. همان‌طور که در بخش ۲-۲-۷ اشاره شد، قدرت پیش‌بینی، اعتبار و تعمیم‌پذیری مدل‌های توسعه یافته در این رساله، با استفاده از پیش‌بینی پاسخ ( $pIC_{50}$ ) داده‌های مجموعه آزمون با مدل بهینه (ALASSO-LM-ANN)، پیش‌بینی پاسخ کل ترکیبات با تکنیک رد مرحله‌ای تک تک (LOO)، محاسبه پارامترهای آماری، دامنه کاربرد و آزمون  $Y$ -تصادفی نیز مورد ارزیابی بیش‌تر قرار گرفت که در ادامه نتایج روش‌های ارزیابی متفاوت آورده شده است.

## ۲-۳-۷-۱ ارزیابی مدل ALASSO-LM-ANN با استفاده از مجموعه آزمون

عدم حضور داده‌های مجموعه آزمون در هر دو مرحله انتخاب توصیف‌کننده‌های مؤثر به روش ALASSO و مدل‌سازی با ANN، معیار مناسبی برای ارزیابی مدل ALASSO-LM-ANN به حساب می‌آیند. به‌این منظور فعالیت دارویی ( $pIC_{50}$ ) مربوط به ۱۳ ترکیب موجود در مجموعه آزمون با استفاده از مدل ANN در شرایط بهینه ALASSO-LM-ANN با معماری ۱-۲-۹ پیش‌بینی شدند و نتایج در جدول ۲-۱۰ آورده شده است. مقادیر خطای کم اغلب ترکیبات موجود در مجموعه آزمون (جدول ۲-۱۰) نشان‌دهنده قدرت پیش‌بینی مناسب مدل توسعه‌یافته ALASSO-LM-ANN است. شکل ۲-۱۰ از رسم مقادیر پیش‌بینی شده  $pIC_{50}$  به‌وسیله مدل بهینه بر حسب مقادیر واقعی  $pIC_{50}$  به دست آمد. با توجه به این که  $R^2$  مربوط به مجموعه آموزش و ارزیابی برای مدل ALASSO-LM-ANN به ترتیب برابر با ۰/۸۲ و ۰/۸۶ است بنابراین با مقایسه مقادیر مذکور با مقدار  $R^2$  مجموعه آزمون (شکل ۲-۱۱) قدرت پیش‌بینی و تعمیم‌پذیری مدل ALASSO-LM-ANN اثبات می‌شود.

جدول ۱۰-۲ نتایج حاصل از ارزیابی مدل ALASSO-LM-ANN با استفاده از مجموعه آزمون

شماره ترکیب	pIC <sub>50</sub>		درصد خطا
	مقدار واقعی	مقدار پیش‌بینی شده	
۳	۴/۸۷	۴/۷۹	-۱/۶۷
۶	۶/۰۱	۶/۷۸	۱۲/۸
۱۹	۴/۸۲	۵/۰۶	۵/۰۵
۲۳	۴/۶	۴/۴۴	-۳/۳۷
۲۹	۴/۸	۴/۶۷	-۲/۷۲
۵۱	۴/۸۶	۴/۷۵	-۲/۲۵
۵۴	۵/۷۷	۴/۹۳	-۱۴/۶
۶۱	۴/۹۲	۴/۸۶	-۱/۳
۶۳	۴/۸۲	۴/۵	-۶/۶۳
۶۵	۴/۸۲	۴/۵۱	-۶/۴
۷۰	۴/۲۲	۴/۴۱	۴/۴
۷۴	۷/۲۲	۶/۶۷	-۷/۶۲
۸۳	۶/۹۱	۶/۶۶	-۳/۵۷



شکل ۱۱-۲ نمودار تغییرات مقادیر پیش‌بینی شده pIC<sub>50</sub> به‌وسیله مدل ALASSO-LM-ANN در شرایط بهینه در مقابل مقادیر تجربی برای داده‌های مجموعه آزمون

۲-۳-۷-۲ ارزیابی مدل ALASSO-LM-ANN با پیش بینی  $pIC_{50}$  تمام ترکیبات مجموعه داده

### با استفاده از روش رد مرحله‌ای تک تک

در این بخش، قدرت پیش‌بینی مدل ALASSO-LM-ANN، با استفاده از تکنیک رد مرحله‌ای تک تک (LOO) برای پیش‌بینی فعالیت دارویی همه ترکیبات مورد بررسی قرار گرفت. در این تکنیک، هر داده یک‌بار به‌عنوان داده آزمون در نظر گرفته شد و مدل ALASSO-LM-ANN با معماری ۱-۲-۹ با ترکیبات باقی‌مانده در مجموعه داده‌ها آموزش داده شد و فعالیت دارویی داده آزمون کنار گذاشته شده، پیش‌بینی شد. این روش ۹۰ بار برای ۹۰ بازدارنده  $3CL^{pro}$  تکرار شد تا زمانی که فعالیت دارویی همه ترکیبات توسط مدل ALASSO-LM-ANN پیش‌بینی شد و نتایج حاصله در جدول ۲-۱۱ خلاصه شد. برای مطالعه بیشتر قدرت پیش‌بینی مدل، مقادیر پیش‌بینی شده  $pIC_{50}$  همه ترکیبات بر حسب مقادیر تجربی آن‌ها رسم شد (شکل ۲-۱۲). همان‌طور که شکل ۲-۱۲ نشان می‌دهد، مقدار  $Q^2_{LOO}$  از حد قابل قبول ( $Q^2=0/5$ ) بالاتر است. بنابراین این نتیجه حاکی از استحکام مناسب مدل می‌باشد. علاوه بر این، مقادیر باقی‌مانده استاندارد شده با استفاده از رابطه ۱-۱۲ محاسبه شد و نتایج آن در مقابل مقادیر واقعی  $pIC_{50}$  رسم شد (شکل ۲-۱۳). با توجه به نمودار باقی‌مانده‌ها، توزیع تصادفی مقادیر باقی‌مانده حول محور صفر، نشان‌دهنده عدم وجود خطای سیستماتیک در مدل پیشنهادی شبکه عصبی توسعه یافته با استفاده از روش ALASSO به‌عنوان روش انتخاب متغیر (ALASSO-LM-ANN) است.

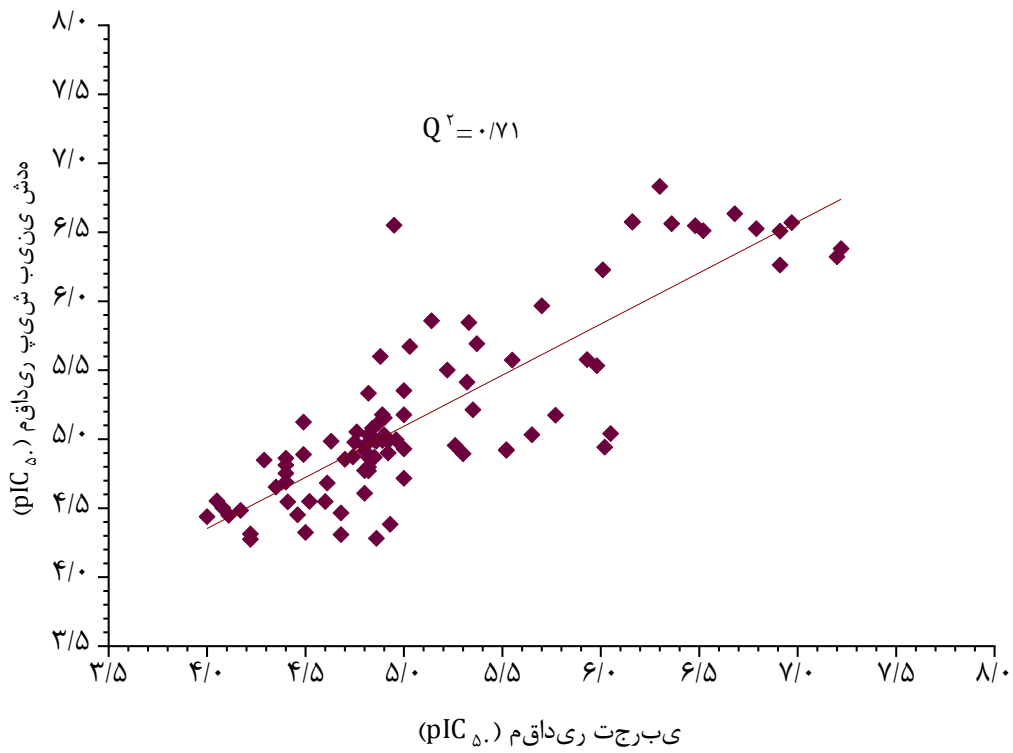


جدول ۱۱-۲ نتایج حاصل از ارزیابی مدل ALASSO-LM-ANN به روش رد مرحله‌ای تک تک برای کل داده‌ها

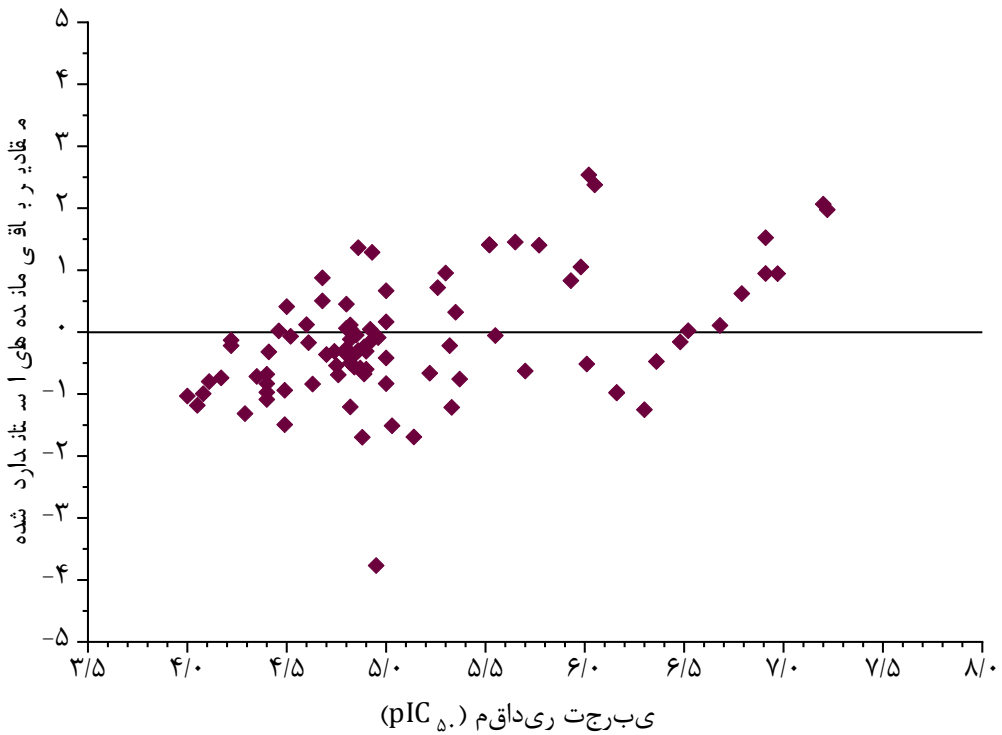
شماره ترکیب	pIC <sub>۵۰</sub>		درصد خطا	شماره ترکیب	pIC <sub>۵۰</sub>		درصد خطا
	مقدار واقعی	مقدار پیش‌بینی شده			مقدار واقعی	مقدار پیش‌بینی شده	
۱	۵/۱۴	۵/۸۶	۱۳/۹۸	۲۹	۴/۸	۴/۹۱	۲/۳۸
۲	۵/۰۳	۵/۶۷	۱۲/۷۶	۳۰	۴/۷۴	۴/۸۷	۲/۷۹
۳	۴/۸۷	۵/۱۲	۵/۰۵	۳۱	۴/۷	۴/۸۵	۳/۲۷
۴	۴/۸۸	۵/۶	۱۴/۷۶	۳۲	۴/۴	۴/۷۵	۸/۰۲
۵	۵/۷	۵/۹۷	۴/۶۸	۳۳	۴/۱۱	۴/۴۵	۸/۲۶
۶	۶/۰۱	۶/۲۳	۳/۶۳	۳۴	۴/۵	۴/۳۲	-۳/۹
۷	۵/۳۲	۵/۴۱	۱/۷۵	۳۵	۴/۴۹	۴/۸۹	۸/۸۸
۸	۶/۰۲	۴/۹۴	-۱۷/۹۱	۳۶	۴/۴۶	۴/۴۵	-۰/۱۹
۹	۴/۹۵	۶/۵۵	۳۲/۳۳	۳۷	۵/۰۰	۵/۳۵	۷/۰۶
۱۰	۴/۶۳	۴/۹۹	۷/۶۸	۳۸	۴/۲۹	۴/۸۵	۱۳/۰۳
۱۱	۴/۹	۵/۱۵	۵/۱۹	۳۹	۵/۳۵	۵/۲۱	-۲/۵۵
۱۲	۴/۷۶	۵/۰۵	۶/۱۷	۴۰	۴/۹	۵/۰۳	۲/۶۷
۱۳	۶/۵۲	۶/۵۱	-۰/۱۴	۴۱	۵/۹۳	۵/۵۸	-۵/۹۴
۱۴	۶/۰۵	۵/۰۴	-۱۶/۶۹	۴۲	۵/۶۵	۵/۰۳	-۱۰/۹۳
۱۵	۵/۲۲	۵/۵	۵/۴	۴۳	۵/۳۷	۵/۶۹	۵/۹۶
۱۶	۴/۹۲	۴/۹	-۰/۳۹	۴۴	۴/۹۳	۴/۳۸	-۱۱/۰۸
۱۷	۴/۸۹	۵/۱۸	۵/۸۷	۴۵	۴/۱۷	۴/۴۸	۷/۵۲
۱۸	۴/۸۹	۴/۹۹	۲/۰۷	۴۶	۴/۰۸	۴/۵	۱۰/۳۴
۱۹	۴/۸۲	۵/۳۳	۱۰/۶۴	۴۷	۴/۸۶	۴/۲۸	-۱۱/۹۱
۲۰	۴/۸	۴/۹۵	۳/۱۴	۴۸	۵/۲۶	۴/۹۶	-۵/۷۹
۲۱	۴/۸	۴/۷۷	-۰/۵۶	۴۹	۴/۸۵	۴/۹۹	۲/۹۴
۲۲	۴/۸	۴/۶۱	۴/۰۰	۵۰	۵/۰۰	۵/۱۸	۳/۵۵
۲۳	۴/۶	۴/۵۵	-۱/۱۲	۵۱	۴/۸۶	۴/۹۹	۲/۶۱
۲۴	۴/۴۹	۵/۱۲	۱۴/۱۳	۵۲	۴/۴	۴/۸۶	۱۰/۵
۲۵	۵/۵۲	۴/۹۲	-۱۰/۸۷	۵۳	۵/۹۸	۵/۵۳	-۷/۴۷
۲۶	۵/۳	۴/۸۹	-۷/۶۴	۵۴	۵/۷۷	۵/۱۷	-۱۰/۳۳
۲۷	۵/۰۰	۴/۹۳	-۱/۴	۵۵	۴/۷۵	۴/۹۸	۴/۸۲
۲۸	۴/۸۲	۴/۸۷	۱/۰۳	۵۶	۵/۵۵	۵/۵۷	۰/۴۲

## ادامه جدول ۱۱-۲

شماره ترکیب	pIC <sub>۵۰</sub>		درصد خطا	شماره ترکیب	pIC <sub>۵۰</sub>		درصد خطا
	مقدار واقعی	مقدار پیش‌بینی شده			مقدار واقعی	مقدار پیش‌بینی شده	
۵۷	۵/۳۳	۵/۸۵	۹/۶۸	۸۴	۶/۹۷	۶/۵۷	-۵/۷۴
۵۸	۵/۵۲	۴/۹۲	-۱۰/۸۱	۸۵	۴/۰۵	۴/۵۵	۱۲/۴
۵۹	۵/۰۰	۴/۷۲	-۵/۶۷	۸۶	۴/۶۱	۴/۶۸	۱/۵۷
۶۰	۴/۹۶	۵/۰۰	۰/۷۵	۸۷	۴/۶۸	۴/۴۷	-۴/۵۷
۶۱	۴/۹۲	۴/۹۹	۱/۳۶	۸۸	۴/۴۱	۴/۵۵	۳/۰۸
۶۲	۴/۸۵	۴/۸۷	۰/۴۱	۸۹	۴/۸۴	۵/۰۸	۴/۹۶
۶۳	۴/۸۲	۴/۷۷	-۱/۰۲	۹۰	۴/۶۸	۴/۳۱	-۷/۹۴
۶۴	۴/۸۲	۵/۰۳	۴/۳۹				
۶۵	۴/۸۲	۴/۸	-۰/۴۴				
۶۶	۴/۵۲	۴/۵۵	۰/۶۵				
۶۷	۴/۴۰	۴/۸۱	۹/۳۸				
۶۸	۴/۴	۴/۶۹	۶/۵۷				
۶۹	۴/۳۵	۴/۶۵	۶/۹۸				
۷۰	۴/۲۲	۴/۲۸	۱/۳۱				
۷۱	۴/۲۲	۴/۳۱	۲/۲۱				
۷۲	۴/۰۰	۴/۴۴	۱۰/۹۶				
۷۳	۷/۲	۶/۳۲	-۱۲/۱۹				
۷۴	۷/۲۲	۶/۳۸	-۱۱/۶۱				
۷۵	۶/۹۱	۶/۵۱	-۵/۸۲				
۷۶	۶/۶۸	۶/۶۳	-۰/۶۸				
۷۷	۶/۳۰	۶/۸۳	۸/۴۴				
۷۸	۶/۷۹	۶/۵۳	-۳/۸۹				
۷۹	۶/۱۶	۶/۵۸	۶/۷۷				
۸۰	۶/۳۶	۶/۵۶	۳/۱۷				
۸۱	۶/۴۸	۶/۵۵	۱/۰۴				
۸۲	۶/۱۶	۶/۵۷	۶/۷۱				
۸۳	۶/۹۱	۶/۲۶	-۹/۳۷				



شکل ۲-۱۲ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌ها بر اساس تکنیک LOO در مقابل مقادیر تجربی



شکل ۲-۱۳ نمودار باقی‌مانده‌های حاصل از پیش‌بینی فعالیت دارویی ترکیبات با استفاده از تکنیک LOO و مقادیر تجربی برحسب مقادیر تجربی

همان‌طور که در بخش ۲-۳-۶ اشاره شد، به منظور مقایسه عملکرد روش انتخاب متغیر ALASSO از روش SR استفاده شد. پس از مقایسه مدل SR-LM-ANN با مدل ALASSO-LM-ANN مشاهده می‌شود که  $R^2_{test}$  و MSE برای مجموعه آزمون مربوط به مدل SR-LM-ANN به ترتیب برابر با ۰/۵۳ و ۰/۴۰ به دست آمد. علاوه بر این  $Q^2_{LOO}$  و MSE مربوط به کل داده‌های پیش‌بینی شده با تکنیک LOO به ترتیب برابر با ۰/۶۷ و ۰/۲۲ به دست آمد. نتایج به دست آمده تأیید می‌کند که مدل ALASSO-LM-ANN دارای قابلیت پیش‌بینی قابل توجهی در مقایسه با مدل SR-LM-ANN است، که این موضوع نشان‌دهنده عملکرد بالای ALASSO در انتخاب توصیف‌کننده‌های مهم در مطالعات QSAR مبتنی بر ANN است

## ۲-۳-۷-۳ ارزیابی مدل ALASSO-LM-ANN با استفاده از پارامترهای آماری

برای بررسی‌های بیش‌تر قدرت پیش‌بینی مدل پیشنهادی ALASSO-LM-ANN، پارامترهای آماری ذکر شده در بخش ۱-۵-۸-۴ برای فعالیت دارویی پیش‌بینی شده ترکیبات مجموعه آزمون و فعالیت دارویی پیش‌بینی شده برای کل ترکیبات به روش رد مرحله‌ای تک تک، محاسبه و در جدول ۲-۱۲ خلاصه شدند. نتایج حاصله (جدول ۲-۱۲) نشان می‌دهد که پارامترهای آماری برای مدل ALASSO-LM-ANN دارای مقادیری در محدوده قابل قبول هستند. به طوری که پارامترهای مربوط به آزمون تروپشا و روی همگی از مقدار هشدار ۰/۵ بزرگ‌تر و به مقدار  $R^2$  نزدیک هستند. شیب نمودار حاصل از مقادیر پیش‌بینی شده بر حسب مقادیر تجربی (و بالعکس) در عرض از مبدأ صفر نیز در محدوده ۰/۸۵ تا ۱/۱۵ قرار دارند که این نتیجه نیز حاکی از صحت مدل توسعه یافته ANN با توصیف‌کننده‌های منتخب روش ALASSO می‌باشد.

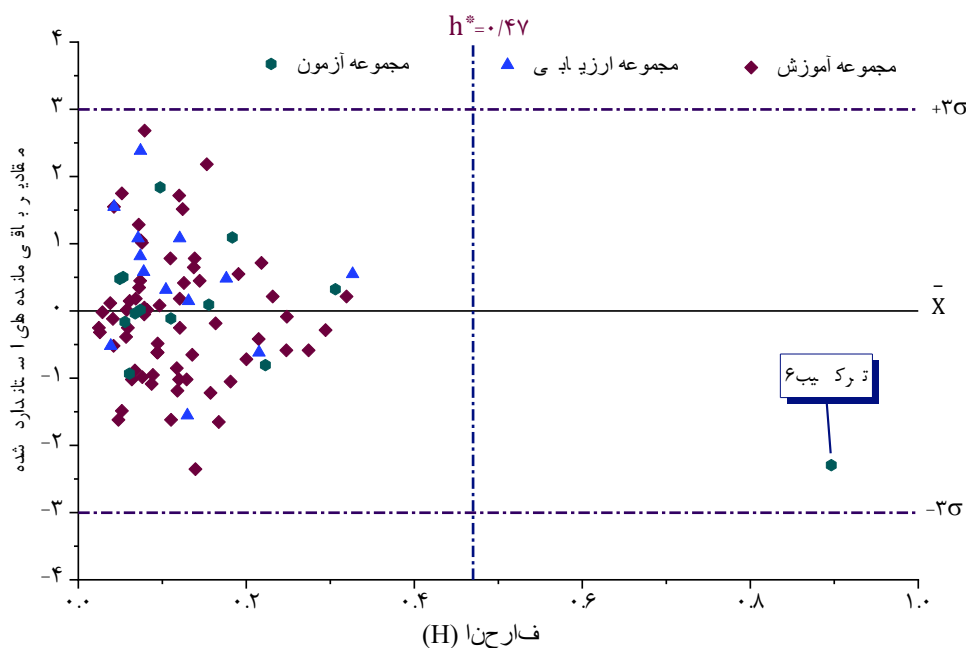
جدول ۲-۱۲ پارامترهای آماری محاسبه شده برای مجموعه آزمون و داده‌های پیش‌بینی شده با تکنیک LOO برای مدل ALASSO-LM-ANN

ردیف	پارامتر آماری	نتایج pIC50 پیش‌بینی شده با ALASSO-LM-ANN		
		ترکیبات مجموعه آزمون	کل ترکیبات به روش LOO	محدوده قابل قبول
۱	PRESS	۱/۴۷	۶/۳۳	-
۲	SEP	۰/۳۴	۰/۴۳	-
۳	MAE	۰/۲۸	۰/۳۳	-
۴	REP(%)	۶/۰۸	۸/۲۷	-
۵	MSE	۰/۱۱	۰/۱۸	-
۶	MRE	۵/۱۴	۶/۳۲	-
۷	R <sup>2</sup>	۰/۸۳	-	R <sup>2</sup> > ۰/۱۶
۸	Q <sup>2</sup> <sub>LOO</sub>	-	۰/۷۱	Q <sup>2</sup> <sub>LOO</sub> > ۰/۵
۹	R <sup>2</sup> <sub>0</sub>	۰/۸۵	۰/۶۳	نزدیک به R <sup>2</sup>
۱۰	R <sup>2</sup> <sub>0</sub> نسبی	۰/۰۱	۰/۱۱	< ۰/۱
۱۱	R <sup>2</sup> <sub>m</sub>	۰/۷۷	۰/۵۱	> ۰/۵
۱۲	R' <sup>2</sup> <sub>0</sub>	۰/۸۶	۰/۷۱	نزدیک به R <sup>2</sup>
۱۳	R' <sup>2</sup> <sub>0</sub> نسبی	۰/۰۰	۰/۰۰	< ۰/۱
۱۴	R' <sup>2</sup> <sub>m</sub>	۰/۷۷	۰/۶۳	> ۰/۵
۱۵	R-R	۰/۰۱	۰/۰۸	< ۰/۳
۱۶	k	۰/۹۸	۱	۰/۸۵ ≤ k ≤ ۱/۱۵
۱۷	k'	۱/۰۳	۰/۹۹	۰/۸۵ ≤ k' ≤ ۱/۱۵

## ۲-۳-۷-۴ ارزیابی ALASSO-LM-ANN با استفاده از دامنه کاربرد

دامنه کاربرد (AD) یک استراتژی ارزشمند برای تأیید اعتبار مدل‌های QSAR توسعه یافته است. به‌منظور نمایش دامنه کاربرد مدل ALASSO-LM-ANN، نمودار ویلیام (شکل ۲-۱۴) مورد استفاده قرار گرفت. بنابراین مطابق با روش کار بخش ۱-۵-۸-۵ و ۲-۲-۷-۴ ابتدا مقدار انحراف (H) با استفاده از رابطه ۱-۱۳ به‌دست آمد. سپس بر اساس مقادیر پیش‌بینی شده pIC<sub>50</sub> توسط مدل ALASSO-LM-

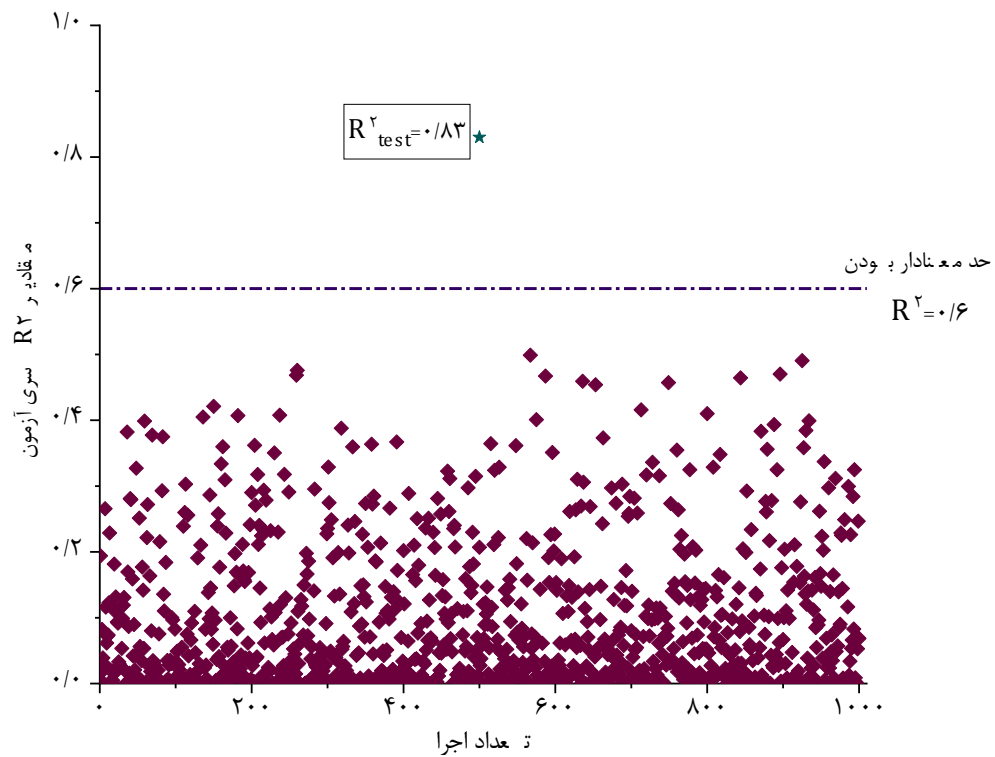
ANN، مقادیر باقی مانده‌ها برای همه ترکیبات محاسبه شد و مقادیر مربوط به باقی مانده‌های استاندارد شده با استفاده از رابطه ۱-۱۴ به دست آمد. نمودار ویلیام از رسم باقی مانده‌های استاندارد شده بر حسب مقادیر انحراف به دست آمد (شکل ۲-۱۴). نمودار ویلیام نشان می‌دهد که مقادیر H محاسبه شده کمتر از حد آستانه‌ای  $h^*$  است. خط عمودی، نشان‌دهنده مقدار آستانه  $h^*$  است که با استفاده از رابطه  $\sqrt{3p/n}$  (برابر با تعداد توصیف‌کننده‌ها به علاوه ۱ (۱۰) و  $n$  نیز برابر با داده‌های مجموعه آموزش و برابر با ۶۴ هست) محاسبه شد و برابر با ۰/۴۷ به دست آمد. علاوه بر این مقادیر باقی مانده‌های استاندارد شده در محدوده  $\pm 3\sigma$  قرار دارند. خطوط افقی، نشان‌دهنده حدود  $\pm 3\sigma$  برای مقادیر باقی مانده‌های استاندارد شده است. نتایج شکل ۲-۱۴ نشان می‌دهد که ۹۹٪ از کل داده‌های به کار رفته در ساخت، ارزیابی و آزمون مدل ALASSO-LM-ANN در فواصل آستانه‌ای دامنه کاربرد قرار دارند. در نتیجه، مجموعه داده‌ها، از جمله مجموعه آموزش، ارزیابی و آزمون در محدوده قابل اعتماد وجود دارند و مدل ALASSO-LM-ANN از اعتبار مناسبی برخوردار است.



شکل ۲-۱۴ دامنه کاربرد مدل ALASSO-LM-ANN، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $\pm 3\sigma$  و  $h^*$  است.

## ۲-۳-۷-۵ ارزیابی مدل ALASSO-LM-ANN با استفاده از آزمون Y-تصادفی

اعتبار ارتباط ساختار-فعالیت دارویی ایجاد شده با مدل برتر ALASSO-LM-ANN و عدم وجود ارتباط تصادفی بین ویژگی‌های ساختار ترکیبات شیمیایی و فعالیت دارویی مربوطه با استفاده از آزمون Y-تصادفی مورد بررسی قرار گرفت. برای انجام آزمون Y-تصادفی، ابتدا مقادیر فعالیت دارویی در محدوده ۷/۲۲-۴/۰۰ به تعداد ۱۰۰۰ بار تصادفی شدند. مدل‌های شبکه عصبی ALASSO-LM-ANN با استفاده از پاسخ‌های تصادفی متغیر وابسته مجموعه آموزش ساخته شدند و برای پیش بینی مقادیر فعالیت دارویی ترکیبات مجموعه آزمون به کار گرفته شدند. مقادیر  $R^2$  حاصل از پیش بینی فعالیت دارویی ترکیبات مجموعه آزمون با ۱۰۰۰ مدل توسعه یافته با متغیر وابسته تصادفی، به دست آمد و بر حسب تعداد اجرا ترسیم شدند (شکل ۲-۱۵). نتایج به دست آمده در شکل ۲-۱۵ نشان می‌دهد که مقادیر  $R^2$  حاصل از هر بار اجرای مدل ALASSO-LM-ANN با داده‌های تصادفی از مقدار قابل قبول ۰/۶ نیز کوچک‌تر هستند [۴۳]. شکل ۲-۱۵ نشان می‌دهد که مقادیر  $R^2$  آزمون Y-تصادفی، در محدوده ۰/۰۰ تا ۰/۴۹ قرار دارند و تنها ۶/۴ درصد از  $R^2$  ها دارای مقادیر بیش‌تر از ۰/۳ هستند. بنابراین مطابق با نتایج شکل ۲-۱۵، مقادیر  $R^2$  داده‌های تصادفی به‌طور قابل توجهی کوچک‌تر از  $R^2$  مجموعه آزمون ( $R^2=0/83$ ) حاصل از مدل پیشنهادی ALASSO-LM-ANN توسعه یافته با متغیر اصلی است. نتایج به دست آمده ثابت می‌کند که ارتباط بین ۹ توصیف کننده منتخب روش ALASSO و فعالیت دارویی ( $pIC_{50}$ ) تصادفی نیست و مدل ALASSO-LM-ANN بر اساس یک رابطه منطقی و معنادار بین متغیرهای مستقل و وابسته ایجاد شده است.



شکل ۲-۱۵ نمودار مقادیر  $R^2$  به دست آمده در آزمون  $Y$ -تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای  $Y$ -تصادفی و پیش‌بینی فعالیت ترکیبات مجموعه آزمون به وسیله مدل ALASSO-LM-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه



## ۲-۴ پیش‌بینی فعالیت دارویی برخی از بازدارنده‌های ایدز و سرطان با

### استفاده از مدل LAD-LASSO-ANN

۲-۴-۱ مقدمه

کمی کردن خواص فیزیکوشیمیایی و فعالیت‌های دارویی ترکیبات دارویی طبیعی و سنتزی و یافتن الگو و روابط حاکم بر کمیت‌های اندازه‌گیری شده همواره موضوع جالب و در حال پیشرفت است که در روش QSAR به آن پرداخته می‌شود. در سال‌های اخیر استفاده از شیمی محاسباتی و شبیه‌سازی مولکولی جهت طراحی دارو به کمک رایانه مورد توجه محققین دارویی قرار گرفته است. از مزایای روش‌های محاسباتی می‌توان به کاهش تعداد ترکیبات ساخته شده برای پیدا کردن ترکیب رهبر، بالا بردن سرعت محاسبات و آزمایش‌ها با پیش‌بینی قابل اعتماد از طریق ویژگی‌های دارویی ساختار مولکول و همچنین کاهش استفاده از آزمایش‌های حیوانات و بالینی اشاره کرد. از این‌رو محققان شیمی، همواره در تلاش هستند تا روش‌های آماری پیشرفته‌تری را برای ایجاد معادلات ریاضی بین داده‌های دارویی و ویژگی‌های ساختاری به کار ببرند. بنابراین ارائه یک مدل QSAR با قابلیت پیش‌بینی و تفسیرپذیری بالا همواره مورد توجه بوده است. یک مدل QSAR با کارایی بالا، باید به ازای استفاده از حداقل متغیرهای وابسته، قدرت پیش‌بینی مناسبی را در مقابل ترکیبات شیمیایی جدید نشان دهد. از این‌رو استفاده از روش‌های انتخاب متغیر کارآمد برای کاهش تعداد توصیف‌کننده‌ها توصیه می‌شود. روش‌های انتخاب متغیر با حذف توصیف‌کننده‌های زائد باعث کاهش توصیف‌کننده‌ها و انتخاب توصیف‌کننده‌های مؤثر می‌شود و تفسیرپذیری مدل افزایش می‌یابد. در طی سال‌های اخیر روش‌های انتخاب متغیر متفاوتی برای گزینش توصیف‌کننده‌هایی با بیش‌ترین ارتباط در مطالعات QSAR استفاده شده است که از این دسته می‌توان به روش‌های انتخاب متغیر کلاسیک و انقباضی اشاره کرد. روش‌های انتخاب متغیر کلاسیک همچون انتخاب پیش‌رو، پس‌رو و روش رگرسیون

گام به گام در مطالعات QSAR مورد توجه بوده است که همه این روش‌ها مبتنی بر روش حداقل مربعات معمولی (OLS) هستند. همان‌طور که واضح است روش OLS در مواجهه با داده‌هایی با تعداد متغیر زیاد و نمونه کم (داده‌هایی با ابعاد بالا و نمونه‌های تنگ) به دلیل وجود همبستگی بین متغیرها یا هم خطی چندگانه موفق عمل نمی‌کنند. علاوه بر این، وجود مشاهدات دور افتاده از جمله عواملی است که برآوردگر OLS را محدود می‌نماید. برای حل این مشکل می‌توان از برآوردگرهای استوار چون برآوردگر کم‌ترین قدر مطلق انحرافات (LAD) بهره گرفت که علاوه بر مقاوم بودن در برابر مشاهدات دور افتاده، برآوردگرهایی با کارایی مطلوب را ایجاد می‌نمایند. بنابراین ونگ<sup>۱</sup> و همکارانش (۲۰۰۷) با افزودن پارامتر جریمه به تابع هدف LAD مدلی را پیشنهاد دادند که علاوه بر مقابله در برابر مشاهدات دور افتاده، قابلیت انتخاب متغیرهای مؤثر را نیز دارد [۳۲]. از جمله روش‌های انتخاب متغیر رگرسیونی انقباضی با تابع جریمه L1 می‌توان به روش انتخاب متغیر انقباضی LASSO [۲۹] اشاره کرد. LASSO به دلیل داشتن تابع جریمه  $\lambda \sum_{i=1}^p |\beta_i|$  ضرایب پارامترهایی بی‌اهمیت را صفر می‌کند و با توجه به این ماهیت ذاتی، توانایی به کارگیری به عنوان یک روش انتخاب متغیر کارآمد را دارد [۲۹]. اخیراً محققین از روش LASSO برای انتخاب توصیف‌کننده‌هایی با بیش‌ترین ارتباط با پاسخ دارویی را در مطالعات QSAR استفاده نموده‌اند و کارایی روش را مورد بررسی قرار داده‌اند [۱۶۴-۱۶۸]. در این مطالعه سعی شده است تا به‌طور هم‌زمان از مزایای ذاتی LASSO و LAD در مواجهه با داده‌های با ابعاد بالا استفاده شود تا مدل، علاوه بر مقابله با مشاهدات دور افتاده بتواند در مواجهه با هم خطی چندگانه نیز به‌طور قابل قبولی عمل نماید.

در راستای تحقیقات پیشین، تلاش شد که با جفت کردن روش انتخاب متغیر LAD-LASSO با روش مدل‌سازی شبکه عصبی، مقادیر فعالیت دارویی برخی از ترکیبات دارویی پیش‌بینی شود. بنابراین در این بخش از مطالعه، سعی شد با به کارگیری یک روش انتخاب متغیر انقباضی جفت شده با شبکه عصبی،

---

<sup>۱</sup>Wang

مدل‌های QSAR با قدرت پیش‌بینی قابل قبول و تفسیرپذیری مناسب ارائه شود. از این روش LAD- LASSO با قابلیت بالایی در انتخاب متغیر بر روی سه نمونه داده‌های دارویی متشکل از بازدارنده‌های ایدز و سرطان کارسینوم کولورکتال انسان و سرطان ریه اجرا شد. توصیف‌کننده‌های منتخب به‌عنوان ورودی روش مدل‌سازی شبکه عصبی مورد استفاده قرار گرفتند. مدل‌های LAD-LASSO-LM-ANN در هر سه نمونه داده از قدرت پیش‌بینی مناسبی برخوردار هستند. بنابراین از مدل‌های توسعه یافته برای پیشنهاد ترکیبات جدید بالقوه استفاده شد.

## ۲-۴-۲ مجموعه داده‌ها

به‌منظور بررسی کارایی روش انتخاب متغیر LAD-LASSO، از سه مجموعه داده استفاده شد. داده‌ها به‌ترتیب متشکل از ۷۳ ترکیب (بازدارنده‌های ایدز) [۱۲۷، ۱۶۹-۱۷۲]، ۷۲ ترکیب (بازدارنده‌های سرطان از سلول کارسینوم کولورکتال انسان) [۸۰، ۸۱، ۱۷۳] و ۷۰ ترکیب (سرطان بافت ریه انسان) هستند [۸۰، ۸۱، ۱۷۳]. فرمت سیستم ورودی خطی ورودی مولکولی ساده شده (SMILES) و فعالیت‌های دارویی مربوط به هر دسته داده در جدول ۲-۱۳ تا جدول ۲-۱۵ خلاصه شده است. فعالیت‌های دارویی مربوط به هر مجموعه داده بر اساس منفی لگاریتم فعالیت دارویی ( $EC_{50}$  برای داده‌های ضد ایدز و  $IC_{50}$  برای هر دو دسته داده ضد سرطان) به‌دست آمد. تقسیم بندی مجموعه داده‌ها با استفاده از الگوریتم KS انجام شد و داده‌ها به سه دسته آموزش، ارزیابی و آزمون تقسیم بندی شدند و در مراحل مختلف مدل‌سازی QSAR مورد استفاده قرار گرفتند.

---

<sup>1</sup>Human colorectal carcinoma

<sup>2</sup>Human lung cancer tissue

جدول ۱۳-۲ مجموعه داده‌های ضد آیدز به همراه مقادیر واقعی و پیش‌بینی شده pEC<sub>50</sub>

ردیف	SMILES	pEC <sub>50</sub>	pEC <sub>50</sub>
		واقعی	پیش‌بینی شده
۱	<chem>c1(nc(nc(c1)NC1CCNCC1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C)C</chem>	۶/۸۵	۶/۶۸
۳۷	<chem>c1(nc(nc(c1)NC1CCNCC1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C#N)C</chem>	۷/۴۲	۶/۸۳
۳۲	<chem>c1(nc(nc(c1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)N)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C)C</chem>	۶/۷۲	۶/۰۹
۴	<chem>c1(nc(nc(c1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)C)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C)C</chem>	۶/۸۹	۷/۰۵
۵	<chem>c1(nc(nc(c1)NC1CCN(CC1)Cc1ccncc1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C)C</chem>	۶/۶۲	۷/۳۲
۶	<chem>c1(nc(nc(c1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)N)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C#N)C</chem>	۶/۸۹	۶/۶۸
۷۷	<chem>c1(nc(nc(c1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)C)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C#N)C</chem>	۷/۲۴	۷/۰۰
۸ <sup>t</sup>	<chem>c1(nc(nc(c1)NC1CCN(CC1)Cc1ccncc1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C#N)C</chem>	۷/۳۳	۷/۲۶
۹	<chem>c1(nc(nc(c1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)N)Nc1ccc(cc1)C#N)Nc1c(cc(cc1C)C)C</chem>	۷/۱۱	۶/۷۷
۱۰	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccncc1)N)Nc1c(cc(cc1C)C)C</chem>	۸/۳۱	۸/۲۳
۱۱ <sup>v</sup>	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)C(=O)OCC)N)Nc1c(cc(cc1C)C)C</chem>	۷/۹۴	۷/۵۷
۱۲	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)C(=O)N)N)Nc1c(cc(cc1C)C)C</chem>	۷/۶۹	۸/۰۳
۱۳	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)N)N)Nc1c(cc(cc1C)C)C</chem>	۸/۲۴	۷/۶۶
۱۴	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)C)N)Nc1c(cc(cc1C)C)C</chem>	۸/۲۶	۷/۷
۱۵	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccccc1)N)Nc1c(cc(cc1C)C)C</chem>	۷/۹۹	۸/۹۳
۱۶	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)[N](=O)O)N)Nc1c(cc(cc1C)C)C</chem>	۸/۱	۷/۹۲
۱۷	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)C#N)N)Nc1c(cc(cc1C)C)C</chem>	۸/۲۵	۸/۴۸
۱۸ <sup>t</sup>	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccncc1)NC)Nc1c(cc(cc1C)C)C</chem>	۸/۱۴	۸/۳
۱۹	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)C(=O)OCC)NC)Nc1c(cc(cc1C)C)C</chem>	۷/۷۸	۷/۸۵
۲۰	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)C(=O)N)NC)Nc1c(cc(cc1C)C)C</chem>	۸/۳۴	۸/۲۵
۲۱ <sup>v</sup>	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)N)NC)Nc1c(cc(cc1C)C)C</chem>	۸/۱۶	۷/۸۹
۲۲ <sup>v</sup>	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)C)NC)Nc1c(cc(cc1C)C)C</chem>	۸/۲۲	۷/۸۸
۲۳	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccncc1)OC)Nc1c(cc(cc1C)C)C</chem>	۷/۸۸	۸/۱۳
۲۴ <sup>t</sup>	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)C(=O)OCC)OC)Nc1c(cc(cc1C)C)C</chem>	۷/۸۳	۷/۹۱
۲۵	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)C(=O)N)OC)Nc1c(cc(cc1C)C)C</chem>	۸/۰۶	۸/۰۱
۲۶	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)N)OC)Nc1c(cc(cc1C)C)C</chem>	۷/۹۷	۷/۹۶
۲۷	<chem>c1(nc(nc(n1)NC1CCN(CC1)Cc1ccc(cc1)S(=O)(=O)C)OC)Nc1c(cc(cc1C)C)C</chem>	۷/۹۲	۷/۸۵
۲۸	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cccc1C)C</chem>	۶/۲۴	۶/۲۸
۲۹	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cc(cc1C)C)C</chem>	۷/۲۵	۶/۲۴
۳۰	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cc(cc1C)C#N)C</chem>	۶/۹۶	۶/۲
۳۱ <sup>t</sup>	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cc(cc1C)Br)C</chem>	۶/۵۷	۶/۴۹
۳۲	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cc(cc1C)Cl)C</chem>	۶/۸	۵/۹۷
۳۳	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cccc1OC)OC</chem>	۵/۳۷	۵/۳۴
۳۴ <sup>v</sup>	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)[N](=O)O)Oc1c(cccc1Cl)Cl</chem>	۶/۲۱	۵/۹۵
۳۵ <sup>t</sup>	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cc(cc1Cl)Cl)Cl</chem>	۶/۷۷	۶/۱
۳۶	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cc(cc1Cl)N(=O)=O)Cl</chem>	۴/۷۷	۵/۰۰
۳۷	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cc(cc1Br)Br)Br</chem>	۶/۹۶	۶/۱۲
۳۸	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cc(cc1Br)C)Br</chem>	۷/۴۷	۶/۶۹
۳۹	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Oc1c(cc(cc1F)F)F</chem>	۶/۱۴	۶/۴۷
۴۰	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Nc1c(cc(cc1C)C)C</chem>	۴/۵۲	۵/۲۶
۴۱	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)N(=O)=O)Nc1c(cc(cc1Br)C)Br</chem>	۵/۴۱	۵/۹۳
۴۲	<chem>c1(c(nc(c1)Nc1ccc(cc1)C#N)[N](=O)O)Nc1c(cc(cc1Br)Br)F</chem>	۵/۵۲	۵/۸۳
۴۳	<chem>c1(c(nc(c1)Nc1ccc(cc1)N(=O)=O)N(=O)=O)Oc1c(cc(cc1C)C#N)C</chem>	۶/۱۴	۶/۰۳

ادامه جدول ۱۳-۲

ردیف	SMILES	pEC <sub>50</sub>	pEC <sub>50</sub>
		واقعی	پیش‌بینی شده
۴۴ <sup>t</sup>	<chem>c1(c(nnc(c1)Nc1ccc(cc1)Cl)Cl)Oc1c(cc(cc1C)C)C</chem>	۶/۸۶	۶/۳۲
۴۵ <sup>v</sup>	<chem>c1(c(nnc(c1)Nc1ccc(cc1)C)Cl)Oc1c(cc(cc1C)C)C</chem>	۶/۷	۶/۲۶
۴۶	<chem>c1(c(nnc(c1)Nc1ccc(cc1)C#N)Cl)Oc1c(cc(cc1C)C)C</chem>	۷/۱	۶/۳۴
۴۷	<chem>c1(c(nnc(c1)Nc1ccc(cc1)C#N)Cl)Oc1c(cc(cc1Cl)Cl)Cl</chem>	۶/۸۱	۵/۹۳
۴۸	<chem>c1(c(nnc(c1)Nc1ccc(cc1)C#N)Cl)Oc1c(cc(cc1Br)Br)Br</chem>	۶/۷۴	۶/۰۸
۴۹	<chem>c1(c(nnc(c1)Nc1ccc(cc1)C#N)Cl)Oc1c(cc(cc1Br)C)Br</chem>	۷/۴۷	۶/۶۳
۵۰	<chem>c1(c(nnc(c1)Nc1ccc(cc1)C#N)Cl)Oc1c(cc(cc1C)Br)C</chem>	۶/۴	۶/۶۰
۵۱	<chem>c1(c(nnc(c1)Nc1ccc(cc1)C#N)Cl)Oc1c(cccc1Cl)Cl</chem>	۶/۶۶	۵/۹۹
۵۲	<chem>c1(c(nnc(c1)Nc1ccc(cc1)C#N)Cl)Oc1c(cc(cc1C)C#N)C</chem>	۶/۸۴	۶/۶۶
۵۳	<chem>c1(c(nnc(c1)Nc1ccc(cc1)C#N)Cl)Oc1c(cccc1C)C</chem>	۶/۲۷	۶/۵۲
۵۴	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)Cl)Oc1c(cc(cc1C)C)C</chem>	۶/۱۱	۵/۹۶
۵۵ <sup>v</sup>	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)C)Oc1c(cc(cc1C)C)C</chem>	۵/۸۱	۶/۱۹
۵۶	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)N(=O)=O)Oc1c(cc(cc1C)C)C</chem>	۵/۲۹	۵/۵۹
۵۷	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)OC)Oc1c(cc(cc1C)C)C</chem>	۵/۸	۶/۴۴
۵۸ <sup>t</sup>	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C)C</chem>	۶/۵۹	۵/۹۳
۵۹ <sup>v</sup>	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1Br)Br)Br</chem>	۶/۲۰	۶/۰۰
۶۰	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1Br)C)Br</chem>	۶/۶۸	۶/۲۶
۶۱ <sup>v</sup>	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)Br)C</chem>	۵/۹۲	۶/۳
۶۲	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)Cl)Oc1c(cc(cc1C)C)C</chem>	۴/۸۴	۴/۹۳
۶۳ <sup>t</sup>	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)C#N)Oc1c(cccc1C)C</chem>	۵/۷۲	۵/۹۵
۶۴	<chem>c1(c(=O)[nH]nc(c1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C#N)C</chem>	۵/۵۵	۵/۴۵
۶۵	<chem>c1(c(=O)n(nc(c1)Nc1ccc(cc1)C#N)C)Oc1c(cc(cc1C)C)C</chem>	۵/۹۴	۵/۹۵
۶۶	<chem>c1(c(=O)n(nc(c1)Nc1ccc(cc1)C#N)C)Oc1c(cc(cc1C)C)C</chem>	۵/۴۷	۶/۰۲
۶۷	<chem>c1(cc(nc(c1)Nc1ccc(cc1)C#N)OC)Oc1c(cc(cc1C)C)C</chem>	۶/۰۸	۶/۱۳
۶۸	<chem>c1(cc(nc(c1)Nc1ccc(cc1)C#N)OC)Oc1c(cc(cc1C)C#N)C</chem>	۶/۱۵	۶/۰۰
۶۹ <sup>t</sup>	<chem>c1(cc(=O)[nH]c(c1)Nc1ccc(cc1)[N](=O)O)Oc1c(cc(cc1C)C)C</chem>	۵/۸۵	۵/۷۵
۷۰	<chem>c1(cc(=O)[nH]c(c1)Nc1ccc(cc1)C)Oc1c(cc(cc1C)C)C</chem>	۵/۸۵	۶/۵۶
۷۱	<chem>c1(cc(=O)[nH]c(c1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C)C</chem>	۶/۴۳	۶/۰۸
۷۲	<chem>c1(cc(=O)[nH]c(c1)Nc1ccc(cc1)C#N)Oc1c(cc(cc1C)C#N)C</chem>	۶/۸۲	۵/۹۴
۷۳ <sup>v</sup>	<chem>c1(cc(=O)[nH]c(c1)Nc1ccc(cc1)C#N)Oc1c(cccc1C)C</chem>	۵/۸۵	۶/۰۲

t و v به ترتیب نمایانگر داده‌های مجموعه آزمون و ارزیابی می‌باشد. سایر داده‌ها مربوط به مجموعه آموزش می‌باشد.

جدول ۲-۱۴ مجموعه داده‌های ضد سرطان کارسینوم کولورکتال به همراه مقادیر واقعی و پیش‌بینی شده pIC<sub>50</sub>

ردیف	SMILES	pIC <sub>50</sub> واقعی	pIC <sub>50</sub> پیش‌بینی شده
۱	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cccc1</chem>	۵/۶۷	۵/۶۲
۲ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cc(ccc1)C</chem>	۶/۳۶	۶/۰۳
۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1ccc(cc1)C</chem>	۶/۳۲	۵/۸۵
۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1ccc(cc1)OC</chem>	۶/۲۴	۵/۶۱
۵ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1c(cccc1)F</chem>	۵/۴۷	۵/۴۶
۶	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cc(ccc1)F</chem>	۵/۶۳	۶/۰۰
۷ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1c(cccc1)Cl</chem>	۵/۸۷	۵/۵۶
۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cc(ccc1)Cl</chem>	۵/۳۶	۵/۷۳
۹ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cc(ccc1)Br</chem>	۵/۳۶	۴/۹۰
۱۰ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1ccc(cc1)Br</chem>	۵/۲۰	۵/۰۰
۱۱ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1c(cccc1)F</chem>	۵/۶۳	۵/۳۵
۱۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cc(c(cc1)Cl)Cl</chem>	۷/۰۹	۶/۷۵
۱۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cc(c(cc1)F)Cl</chem>	۶/۷۲	۷/۱۳
۱۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C</chem>	۴/۲۳	۴/۵۲
۱۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC</chem>	۴/۲۸	۴/۲۵
۱۶ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CCC</chem>	۴/۲۱	۴/۲۳
۱۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C(C)C</chem>	۴/۱۸	۴/۴۹
۱۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CCCC</chem>	۴/۳۷	۴/۲۲
۱۹ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(C)C</chem>	۴/۲۰	۴/۴۱
۲۰ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CCCCC</chem>	۴/۴۵	۴/۲۷
۲۱ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C1CCCC1</chem>	۴/۵۴	۴/۶۵
۲۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C1CCCCC1</chem>	۴/۶۵	۴/۵
۲۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC=C</chem>	۴/۳۰	۴/۴۸
۲۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=C)C</chem>	۴/۳۵	۴/۴۴
۲۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC=C(C)C</chem>	۴/۵۴	۴/۷۴
۲۶	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1cccc1</chem>	۴/۲۰	۴/۴۹
۲۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)C</chem>	۴/۱۷	۴/۳۴
۲۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)OC</chem>	۴/۲۱	۴/۱۳
۲۹ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)c1cccc1</chem>	۵/۴۷	۵/۵۹
۳۰	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)F</chem>	۴/۳۰	۴/۳۳
۳۱ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)Br</chem>	۴/۲۸	۴/۱
۳۲ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1cccc1Cl</chem>	۴/۷۴	۴/۹۸
۳۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1cccc(c1)Cl</chem>	۴/۵۹	۴/۲۳
۳۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)Cl</chem>	۴/۷۱	۴/۶۴
۳۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)F</chem>	۵/۱۰	۵/۳۲
۳۶	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)Cl</chem>	۴/۴۹	۴/۶۷
۳۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(c(c1)Cl)Cl</chem>	۵/۳۴	۵/۵۵

ادامه جدول ۱۴-۲

ردیف	SMILES	pIC <sub>50</sub> واقعی	pIC <sub>50</sub> پیش‌بینی شده
۳۸ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1cccc(c1)N(=O)=O</chem>	۴/۷۶	۵/۱۲
۳۹	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)C(F)(F)F</chem>	۴/۷۴	۴/۹۳
۴۰ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1cccc1</chem>	۴/۶۶	۴/۵۶
۴۱	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)C</chem>	۴/۷۴	۴/۶۴
۴۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)OC</chem>	۴/۷۲	۵/۰۶
۴۳ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)c1cccc1</chem>	۶/۶۷	۶/۷۷
۴۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)F</chem>	۵/۷۶	۵/۶۲
۴۵ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)Br</chem>	۴/۷۹	۴/۹۲
۴۶	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1cccc1Cl</chem>	۵/۲۹	۵/۵
۴۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1cccc(c1)Cl</chem>	۵/۳۹	۵/۱۴
۴۸ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)Cl</chem>	۵/۶۵	۵/۴۲
۴۹	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)F</chem>	۶/۴۴	۶/۲۹
۵۰	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)Cl</chem>	۵/۶۵	۵/۴۷
۵۱ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(c(c1)Cl)Cl</chem>	۷/۰۰	۶/۸
۵۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1cccc(c1)N(=O)=O</chem>	۶/۲۰	۶/۳۱
۵۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)C(F)(F)F</chem>	۶/۰۰	۶/۱۷
۵۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1cccc1</chem>	۴/۳۷	۴/۴۹
۵۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(C)cc1</chem>	۵/۰۹	۵/۳۲
۵۶ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)OC</chem>	۵/۲	۴/۸۶
۵۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)C(C)(C)C</chem>	۵/۱۷	۵/۲۷
۵۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)F</chem>	۴/۴۸	۴/۷۱
۵۹	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)Cl</chem>	۴/۵۶	۵/۰۵
۶۰ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)Br</chem>	۴/۷۱	۴/۴۳
۶۱	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1c(ccc1)N(=O)=O</chem>	۴/۹۳	۵/۱۴
۶۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1cc(ccc1)N(=O)=O</chem>	۵/۳۸	۵/۸۴
۶۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)N(=O)=O</chem>	۵/۹۸	۵/۷۹
۶۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)C(F)(F)F</chem>	۵/۵۷	۴/۹۹
۶۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1c(ccc(c1)C)C</chem>	۶/۴۴	۶/۴۲
۶۶	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1cc(c(cc1)Cl)N(=O)=O</chem>	۶/۹۲	۶/۸۱
۶۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1c(C)cc(C)cc1C</chem>	۷/۰۹	۶/۷
۶۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1c(cc(cc1C(C)C)C(C)C)C(C)C</chem>	۷/۱۵	۷/۳
۶۹ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)CC</chem>	۴/۰۸	۴/۱۵
۷۰	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)CCC</chem>	۴/۱۷	۴/۱۳
۷۱	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)C1CC1</chem>	۴/۰۱	۴/۴۷
۷۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)CCCCl</chem>	۴/۰۷	۴/۲۸

t و v به ترتیب نمایانگر داده‌های مجموعه آزمون و ارزیابی را می‌باشد. سایر داده‌ها مربوط به مجموعه آموزش می‌باشد.

جدول ۱۵-۲ مجموعه داده‌های ضد سرطان ریه به همراه مقادیر واقعی و پیش‌بینی شده pIC<sub>50</sub>

ردیف	SMILES	pIC <sub>50</sub> واقعی	pIC <sub>50</sub> پیش‌بینی شده
۱ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cccc1</chem>	۵/۱۴	۴/۸۷
۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1ccc(ccc1)C</chem>	۶/۰۷	۶/۱۴
۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1ccc(cc1)C</chem>	۶/۱۱	۵/۷۲
۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1ccc(cc1)OC</chem>	۵/۷۵	۵/۰۴
۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1c(cccc1)F</chem>	۴/۹۹	۵/۴۹
۶	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1ccc(ccc1)F</chem>	۵/۶۵	۶/۱۸
۷ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1c(cccc1)Cl</chem>	۵/۴۸	۴/۸۵
۸ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1ccc(ccc1)Cl</chem>	۵/۰۳	۵/۱۲
۹	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cc(ccc1)Br</chem>	۵/۲	۴/۹۴
۱۰	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1ccc(cc1)Br</chem>	۴/۹۱	۵/۱
۱۱ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1c(cccc1)F</chem>	۵/۳۷	۵/۵۴
۱۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cc(c(cc1)Cl)Cl</chem>	۶/۵۵	۶/۲۴
۱۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)Cc1cc(c(cc1)F)Cl</chem>	۶/۵۱	۶/۴۶
۱۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C</chem>	۴/۰۴	۴/۱۷
۱۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC</chem>	۴/۰۷	۴/۱۳
۱۶ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CCC</chem>	۴/۰۵	۴/۰۹
۱۷ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C(C)C</chem>	۴/۱۷	۴/۵۷
۱۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CCCC</chem>	۴/۱۴	۳/۹۷
۱۹ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(C)C</chem>	۳/۹۴	۴/۰۰
۲۰	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CCCCC</chem>	۴/۲۳	۴/۰۱
۲۱	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C1CCCC1</chem>	۴/۶۱	۴/۶۱
۲۲ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C1CCCCC1</chem>	۴/۷۱	۴/۳۲
۲۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC=C</chem>	۴/۱۱	۴/۳۴
۲۴ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=C)C</chem>	۴/۴۶	۴/۳
۲۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC=C(C)C</chem>	۴/۵	۴/۹۵
۲۶ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1cccc1</chem>	۴/۰۱	۴/۳۵
۲۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)C</chem>	۴/۱	۴/۱۱
۲۸ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)OC</chem>	۴/۱۶	۳/۷۵
۲۹	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)c1cccc1</chem>	۵/۱۸	۵/۲۴
۳۰	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)F</chem>	۴/۰۷	۴/۰۳
۳۱	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)Br</chem>	۴/۲۱	۴/۶۲
۳۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1cccc1Cl</chem>	۴/۷۵	۴/۶۶
۳۳ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(c1)Cl</chem>	۴/۴۹	۴/۵۸
۳۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)Cl</chem>	۴/۳۴	۴/۶۷
۳۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)F</chem>	۵/۰۳	۴/۹۹
۳۶ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)Cl</chem>	۴/۵۶	۴/۵۸
۳۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(c(c1)Cl)Cl</chem>	۴/۸۸	۴/۵۶
۳۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1cccc(c1)N(=O)=O</chem>	۴/۴۵	۴/۲۲
۳۹	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)CC(=O)c1ccc(cc1)C(F)(F)F</chem>	۴/۵	۴/۹۱
۴۰	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1cccc1</chem>	۴/۴۷	۴/۵۸
۴۱	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)C</chem>	۴/۵۵	۴/۷۸
۴۲ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)OC</chem>	۴/۴۱	۴/۸۷
۴۳ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)c1cccc1</chem>	۶/۲۱	۶/۱۵
۴۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)F</chem>	۵/۶۱	۵/۶۴
۴۵ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)Br</chem>	۴/۴۳	۴/۷۷
۴۶	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1cccc1Cl</chem>	۵/۰۷	۴/۶۱
۴۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(c1)Cl</chem>	۴/۹۸	۵/۳۷
۴۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)Cl</chem>	۵/۲	۴/۶۶
۴۹	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)F</chem>	۶/۱۱	۵/۶۳



ردیف	SMILES	pIC <sub>50</sub> واقعی	pIC <sub>50</sub> پیش‌بینی شده
۵۰	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1Cl)Cl</chem>	۵/۲	۴/۸
۵۱	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(c(c1)Cl)Cl</chem>	۶/۵۱	۵/۹۳
۵۲	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1cccc(c1)N(=O)=O</chem>	۵/۶۵	۶/۰۱
۵۳ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)C/C(=N/O)/c1ccc(cc1)C(F)(F)F</chem>	۵/۴۹	۶/۰۹
۵۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccccc1</chem>	۴/۲۲	۴/۶
۵۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(C)cc1</chem>	۵/۰۱	۴/۶۶
۵۶	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)OC</chem>	۵/۰۳	۴/۵۸
۵۷	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)C(C)(C)C</chem>	۵/۱۱	۵/۶
۵۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)F</chem>	۴/۳	۴/۷۴
۵۹	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)Cl</chem>	۴/۴۹	۴/۸۵
۶۰ <sup>v</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)Br</chem>	۴/۶۶	۴/۵۸
۶۱	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1c(ccc1)N(=O)=O</chem>	۴/۸۶	۴/۶۸
۶۲ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1cc(ccc1)N(=O)=O</chem>	۵/۳۵	۵/۲
۶۳	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)N(=O)=O</chem>	۶/۰۱	۵/۹
۶۴	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1ccc(cc1)C(F)(F)F</chem>	۵/۷۲	۵/۴۴
۶۵	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1c(ccc(c1)C)C</chem>	۶/۴۶	۶/۰۲
۶۶ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1cc(c(cc1)Cl)N(=O)=O</chem>	۶/۳۲	۵/۹۳
۶۷ <sup>t</sup>	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1c(C)cc(C)cc1C</chem>	۶/۷۷	۶/۳
۶۸	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)c1c(cc(cc1C(C)C)C(C)C)C(C)C</chem>	۶/۹۲	۶/۱۶
۶۹	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)CC</chem>	۴/۰۵	۴/۵۲
۷۰	<chem>n1(nc2c(c1)c(=O)oc1c2cccc1)S(=O)(=O)CCC</chem>	۴/۱	۴/۵۴

t و v به ترتیب نمایانگر داده‌های مجموعه آزمون و ارزیابی می‌باشد. سایر داده‌ها مربوط به مجموعه آموزش می‌باشد.

## ۲-۴-۳ رسم و بهینه‌سازی ساختار ترکیبات شیمیایی مجموعه داده‌های متفاوت

برای محاسبه توصیف‌کننده‌های مولکولی با مقادیر صحیح، ساختار ترکیبات مورد مطالعه با استفاده

از نرم‌افزار هایپرکم و مطابق با روش کار بخش ۱-۵-۳ بهینه شدند. ساختارهای بهینه شده با پسوند \*.hin\*

ذخیره شدند و به‌عنوان ورودی نرم‌افزار محاسباتی دراگون استفاده شدند.

## ۲-۴-۴ استخراج توصیف‌کننده‌های ساختاری

ساختارهای بهینه شده ترکیبات مورد مطالعه برای هر سه مجموعه داده‌ها به‌طور مجزا در نرم‌افزار

دراگون فراخوانی شدند و سپس برای هر ترکیب تعداد ۳۲۲۴ توصیف‌کننده در ۲۲ دسته متفاوت محاسبه

شدند.

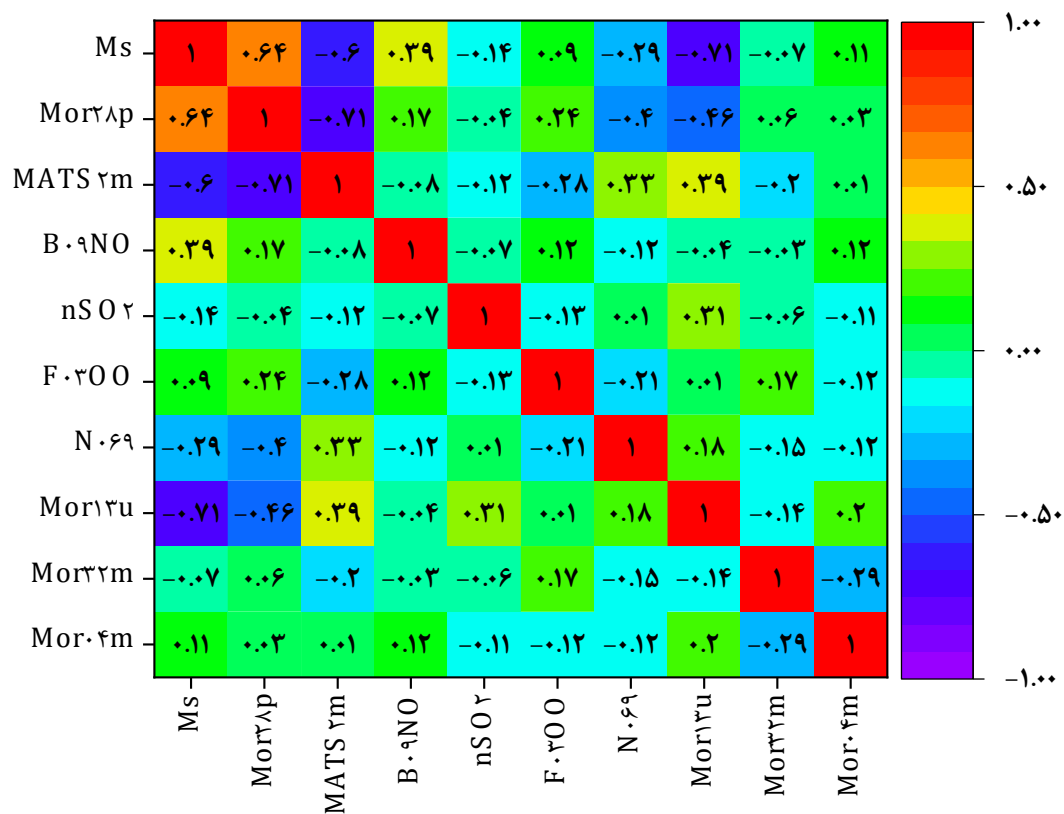
## ۲-۴-۵ پیش‌پردازش و انتخاب توصیف‌کننده‌های مؤثر

با توجه به این‌که وجود توصیف‌کننده‌هایی با مقادیر ثابت و نسبتاً ثابت اطلاعات مفیدی را به مدل نمی‌افزاید، بنابراین این نوع از توصیف‌کننده‌ها (توصیف‌کننده‌هایی با واریانس کمتر از ۰/۰۰۱) با استفاده از بسته نرم‌افزاری caret در نرم‌افزار R حذف شدند. علاوه بر این توصیف‌کننده‌هایی با همبستگی بالای ۰/۹ نیز در نرم‌افزار متلب مورد بررسی قرار گرفتند و از بین دو توصیف‌کننده با همبستگی بالای ۰/۹، توصیف‌کننده‌ای با بیش‌ترین همبستگی با پاسخ، نگه داشته شد و توصیف‌کننده بعدی حذف شد. پس از انجام مراحل پیش‌پردازش داده‌ها، تعداد توصیف‌کننده‌ها برای مجموعه داده‌های ۷۳، ۷۲ و ۷۰ تایی به ترتیب از ۳۲۲۴ توصیف‌کننده به ۳۳۵، ۴۳۸ و ۴۲۹ توصیف‌کننده کاهش یافت. این توصیف‌کننده‌ها به‌طور مستقل به‌عنوان ورودی روش انتخاب متغیر LAD-LASSO مورد استفاده قرار گرفتند. به‌طوری‌که توصیف‌کننده‌های حاصل از مرحله پیش‌پردازش به‌عنوان متغیرهای مستقل و فعالیت‌های دارویی مربوطه به‌عنوان متغیر وابسته در روش انتخاب متغیر LAD-LASSO به‌کار گرفته شدند. به‌منظور اجرای روش انتخاب متغیر LAD-LASSO، مجموعه داده‌ها با استفاده از الگوریتم KS به سه بخش آموزش، ارزیابی و آزمون تقسیم شدند (شماره ترکیبات مربوط به هر مجموعه در جدول مجموعه داده‌ها مشخص شده است). در مرحله انتخاب متغیر، داده‌های مجموعه آزمون از ابتدا بیرون گذاشته شدند و از مجموعه داده‌های مجموعه آموزش و ارزیابی برای انتخاب توصیف‌کننده‌های مؤثر، استفاده شد. روش انتخاب متغیر انقباضی با روش ارزیابی تقاطعی ده فولد موجود در بسته quantreg در نرم‌افزار R روی مجموعه داده‌های ارزیابی و آموزش اجرا شد تا موثرترین توصیف‌کننده‌ها انتخاب شود. با اجرای روش LAD-LASSO در نهایت توصیف‌کننده‌هایی با ضرایب غیر صفر (۱۰، ۱۴ و ۹ توصیف‌کننده به‌ترتیب برای مجموعه داده‌های ۷۳، ۷۲ و ۷۰ تایی) در  $\lambda_{\min}$  مربوطه استخراج شدند. نام و نوع این توصیف‌کننده‌ها به همراه ضرایب رگرسیونی استاندارد شده مربوط به هر توصیف‌کننده در جدول ۲-۱۶ خلاصه شده‌اند. به‌منظور بررسی آماری دقیق‌تر

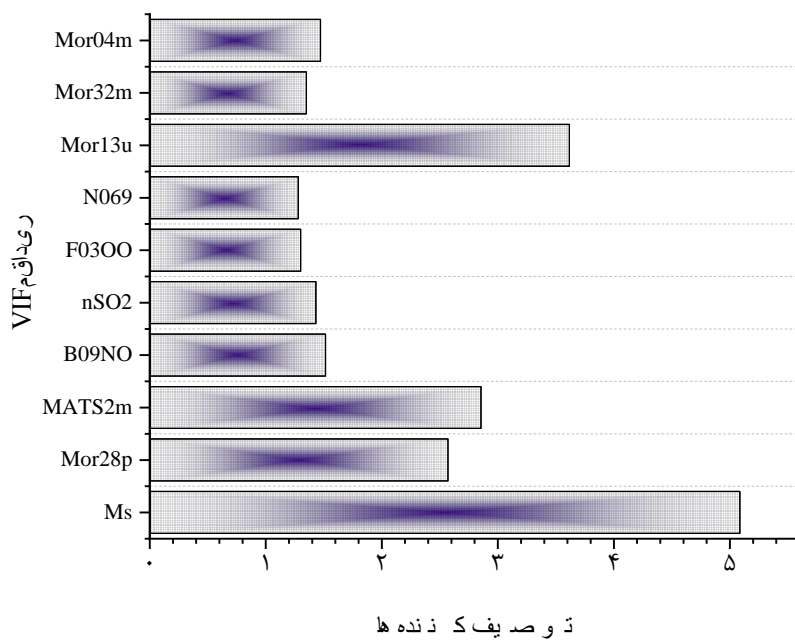
توصیف‌کننده‌های منتخب روش LAD-LASSO هر سه مجموعه داده‌ها، احتمال وجود همبستگی و هم‌خطی بین توصیف‌کننده‌ها به ترتیب با محاسبه مقادیر ضرایب همبستگی بین دو توصیف‌کننده و مقادیر افزایش تورم واریانس (VIF) توصیف‌کننده (مطابق با رابطه ۱-۱۰) مطالعه شد. به این منظور نمودارهای نقشه رنگی و VIF برای توصیف‌کننده‌های منتخب روش LAD-LASSO برای هر سه مجموعه داده رسم شد و نتایج حاصل در شکل ۲-۱۶ تا شکل ۲-۲۱ نمایش داده شده‌اند. نتایج نمودار نقشه رنگی (شکل ۲-۱۶ و شکل ۲-۱۸ و شکل ۲-۲۰) نشان می‌دهد که بین توصیف‌کننده‌های منتخب هر سه مجموعه داده‌ها به روش LAD-LASSO همبستگی معناداری وجود ندارد. علاوه بر این، نمودارهای VIF هر سه مجموعه داده‌ها (شکل ۲-۱۷ و شکل ۲-۱۹ و شکل ۲-۲۱) نیز نشان می‌دهد که مقادیر VIF محاسبه شده برای متغیرهای منتخب هر سه مجموعه داده‌ها، در محدوده ۱-۱۰ قرار دارند. این شواهد بیانگر این است که بین توصیف‌کننده‌های منتخب روش LAD-LASSO هم‌خطی شدیدی وجود ندارد.

جدول ۲-۱۶ توصیف‌کننده‌های منتخب LAD-LASSO برای هر سه مجموعه داده‌ها

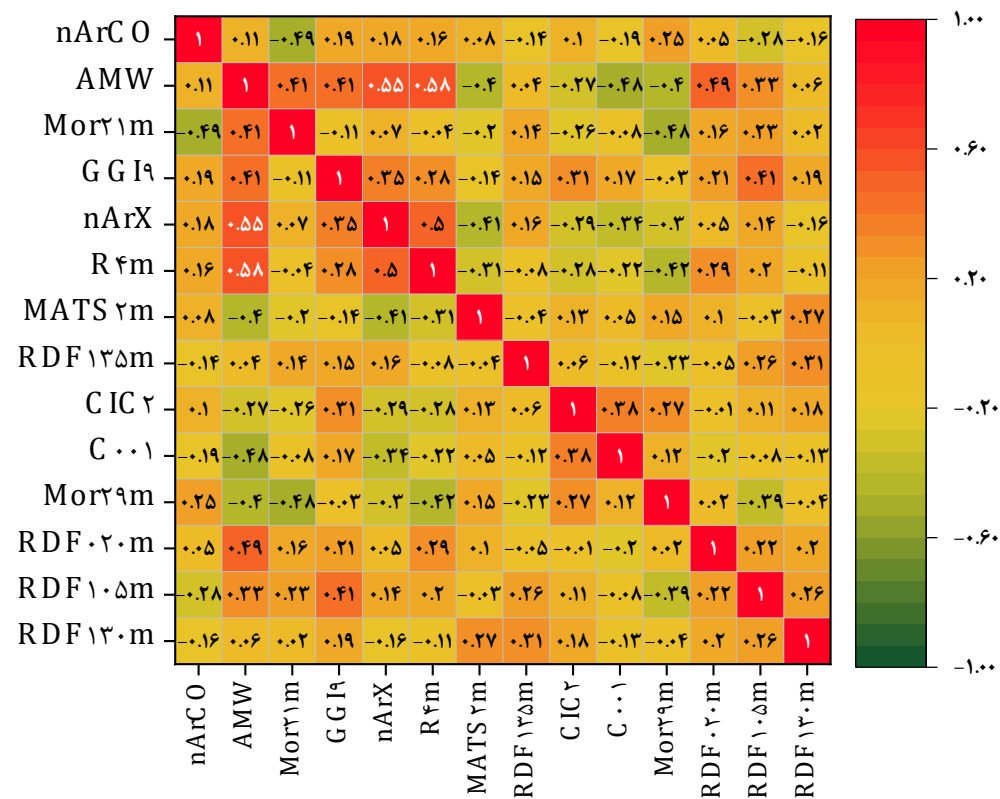
مجموعه داده‌ها	ردیف	نماد	معنا	طبقه‌بندی	ضرایب استاندارد شده
بازدارنده‌های آیدز	۱	Ms	Electro-topological State signal 28 / weighted by polarizability	Constitutional indices	۰/۶۵
	۲	Mor28p	Moran autocorrelation of lag 2 weighted by mass	3D-MoRSE descriptors	-۰/۵۱
	۳	MATS2m	Presence/absence of N - O at topological distance 9	2D autocorrelations	۰/۳۸
	۴	B09NO	number of sulfites (thio-/dithio-)	2D binary fingerprints	-۰/۳۷
	۵	nSO2	Frequency of C - O at topological distance 3	Functional group counts	۰/۲۷
	۶	F03OO	Ar-NH2 / X-NH2	2D frequency fingerprints	-۰/۲۴
	۷	N069	signal 13 / unweighted	Atom-centered fragments	۰/۲۲
	۸	Mor13u	signal 32 / weighted by mass	3D-MoRSE descriptors	۰/۱۹
	۹	Mor32m	signal 04 / weighted by mass	3D-MoRSE descriptors	۰/۱۷
	۱۰	Mor04m	number of ketones (aromatic)	3D-MoRSE descriptors	-۰/۱۴
بازدارنده‌های سرطان کارسینوم کولور کتال	۱	nArCO	average molecular weight	Functional group counts	-۰/۶۵
	۲	AMW	signal 21 / weighted by mass	Constitutional indices	۰/۵۲
	۳	Mor21m	topological charge index of order 9	3D-MoRSE descriptors	-۰/۴۲
	۴	GGI9	number of X on aromatic ring	2D autocorrelations	۰/۳۹
	۵	nArX	R autocorrelation of lag 4 / weighted by mass	Functional group counts	۰/۳۵
	۶	R4m	Moran autocorrelation of lag 2 weighted by mass	GETAWAY descriptors	-۰/۲۷
	۷	MATS2m	Radial Distribution Function - 135 / weighted by mass	2D autocorrelations	۰/۱۷
	۸	RDF135m	Complementary Information Content index (neighborhood symmetry of 2-order)	RDF descriptors	۰/۱۶
	۹	CIC2	CH3R / CH4	Information indices	۰/۱۴
	۱۰	C001	signal 29 / weighted by mass	Atom-centered fragments	۰/۱۴
	۱۱	Mor29m	Radial Distribution Function - 020 / weighted by mass	3D-MoRSE descriptors	۰/۱
	۱۲	RDF020m	Radial Distribution Function - 105 / weighted by mass	RDF descriptors	-۰/۰۹
	۱۳	RDF105m	Radial Distribution Function - 130 / weighted by mass	RDF descriptors	-۰/۰۸
	۱۴	RDF130m	average molecular weight	RDF descriptors	۰/۰۶
بازدارنده‌های سرطان ریه	۱	AMW	Radial Distribution Function - 110 / weighted by mass	Constitutional indices	۰/۶۱
	۲	RDF110m	signal 12 / weighted by Sanderson electronegativity	RDF descriptors	۰/۵۶
	۳	Mor12e	number of terminal primary C(sp3)	3D-MoRSE descriptors	-۰/۵۴
	۴	nCp	signal 29 / weighted by mass	Functional group counts	۰/۴۴
	۵	Mor29m	F attached to C1(sp2)	3D-MoRSE descriptors	۰/۳۳
	۶	F084	Radial Distribution Function - 130 / weighted by mass	Atom-centered fragments	۰/۳۱
	۷	RDF130m	number of ethers (aromatic)	RDF descriptors	۰/۱۵
	۸	nArOR	Radial Distribution Function - 070 / weighted by mass	Functional group counts	-۰/۱۳
	۹	RDF070m		RDF descriptors	-۰/۰۸



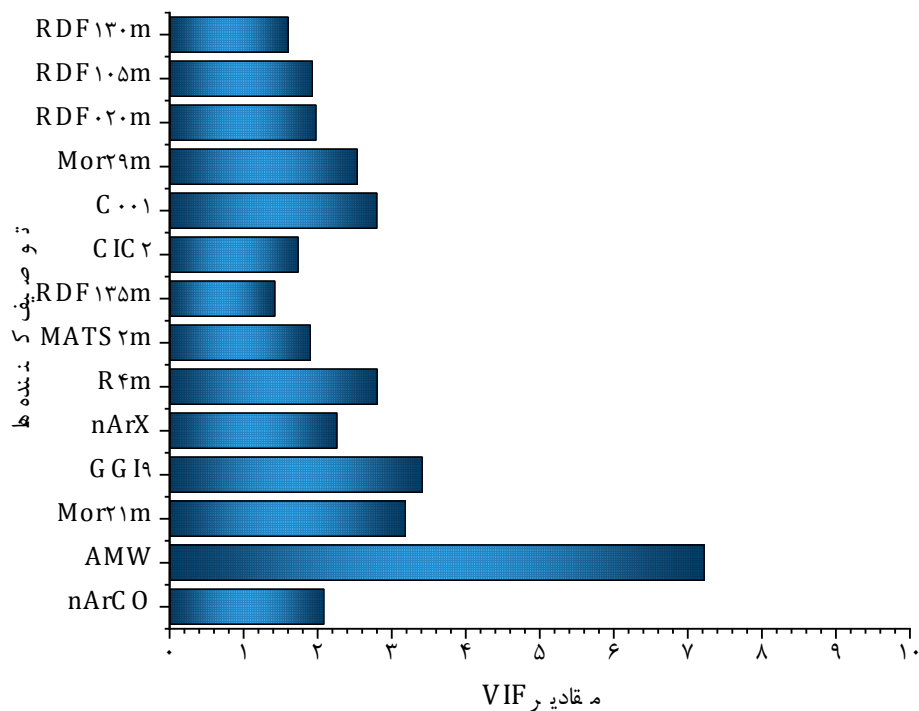
شکل ۲-۱۶ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های ایدز



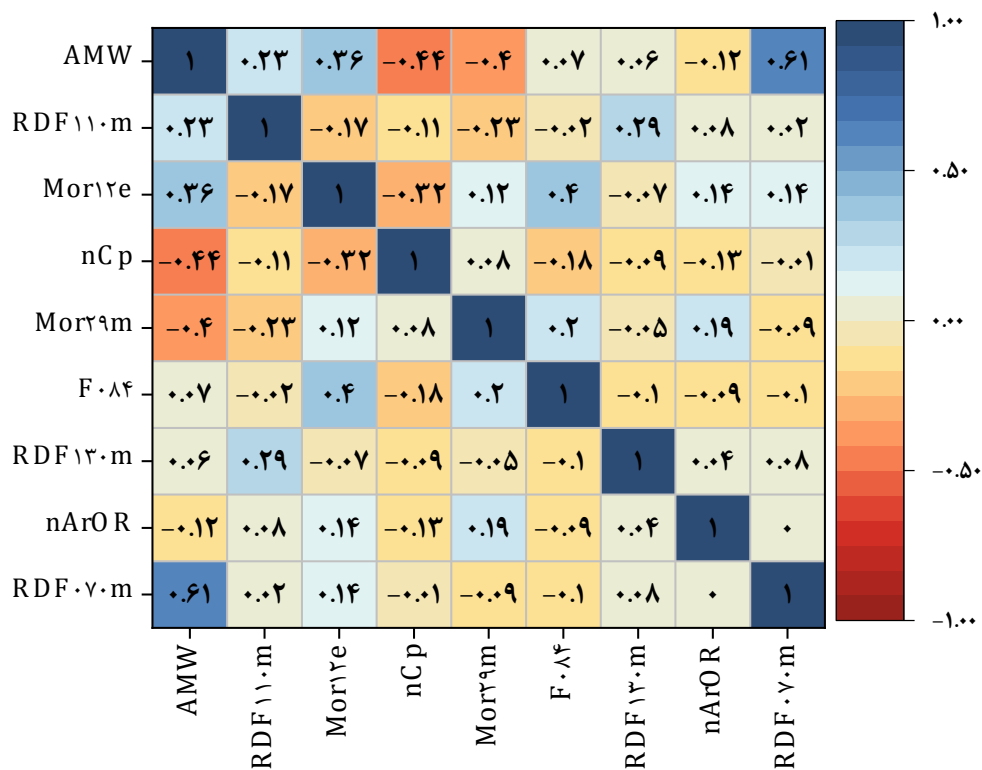
شکل ۲-۱۷ نمودار مقادیر VIF توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های ایدز



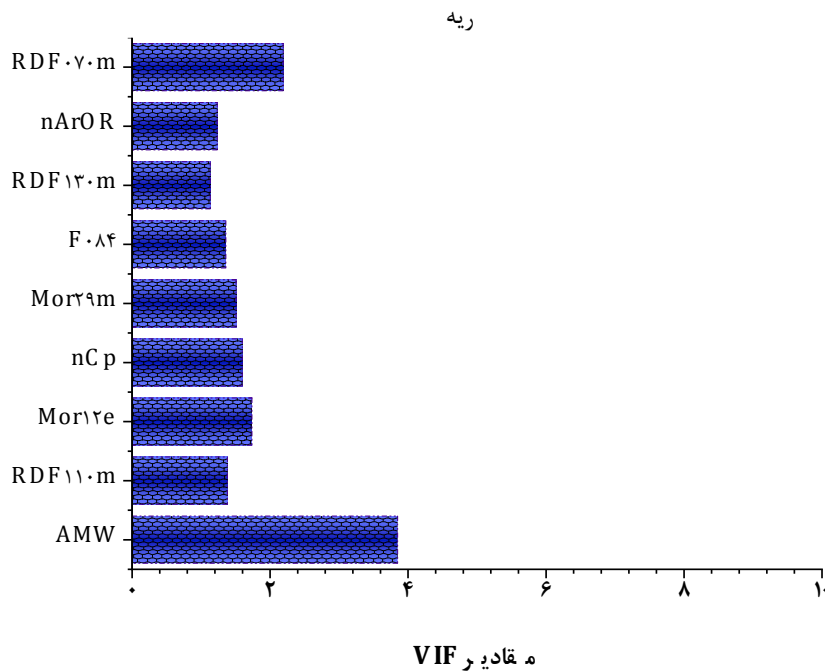
شکل ۲-۱۸ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های سرطان کارسینوم کولورکتال



شکل ۲-۱۹ نمودار مقادیر VIF توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های سرطان کارسینوم کولورکتال



شکل ۲۰-۲ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های سرطان



شکل ۲۱-۲ نمودار مقادیر VIF توصیف‌کننده‌های منتخب LAD-LASSO برای بازدارنده‌های سرطان ریه

## ۲-۴-۶ مدل سازی شبکه عصبی با استفاده از توصیف کننده های منتخب -LAD

### LASSO

برای برقراری ارتباط بین توصیف کننده های منتخب و متغیر وابسته مربوطه از مدل شبکه عصبی مصنوعی سه لایه ای با الگوریتم پس انتشار خطا استفاده شد. به این منظور، یک شبکه پرسپترون، با یک لایه ورودی، یک لایه پنهان و یک لایه خروجی طراحی شد [۱۴۰، ۱۴۱]. برای انتخاب بهترین ساختار شبکه عصبی، بهینه سازی تعداد نورون های لایه ورودی، تعداد نورون های لایه پنهان، تعداد دور آموزش، توابع انتقال (تانژانت هایپربولیک سیگموئیدی و لگاریتم سیگموئیدی که به ترتیب با توابع  $\text{tansig}$  و  $\text{logsig}$  در جعبه ابزار متلب مشخص می شوند) و توابع آموزش (تنظیم بایزین و لونیگ-مارکوارت که به ترتیب با توابع  $\text{trainlm}$  و  $\text{trainbr}$  در جعبه ابزار متلب مشخص می شوند) به طور هم زمان انجام شد. علاوه بر این از تابع خطی  $\text{purlin}$  نیز به عنوان لایه خروجی استفاده شد. بنابراین، توصیف کننده های منتخب بر اساس بزرگی ضرایب استاندارد شده LAD-LASSO (جدول ۲-۱۶) چیده شدند و به عنوان ورودی در ساخت مدل شبکه عصبی مورد استفاده قرار گرفتند. معیار بهینه سازی حداقل نمودن پارامتر  $\text{MSE}$  مجموعه ارزیابی است. از این رو برای بهینه سازی پارامترهای شبکه عصبی مصنوعی، تعداد توصیف کننده ها، گره ها و دورهای آموزشی به طور هم زمان از ۲ تا تعداد کل توصیف کننده های منتخب (با گام ۱)، ۲ تا ۱۰ (با گام ۱) و ۵ تا ۵۰ (با گام ۵) تغییر یافت. ساختارهای شبکه عصبی با ورودی های متفاوت (۳۲۴۰ حالت، ۴۶۸۰ حالت و ۲۸۸۸ حالت برای مجموعه داده های متشکل از ۷۳، ۷۲ و ۷۰ ترکیب دارویی) برای هر سه مجموعه داده ها طراحی شد. آموزش شبکه عصبی با استفاده از داده های مجموعه آموزش و با توابع آموزش تنظیم بایزین و لونیگ-مارکوارت انجام شد. پس از بهینه سازی پارامترهای شبکه عصبی، ساختار بهینه برای هر مجموعه مورد مطالعه، با توجه به حداقل مقدار  $\text{MSE}$  مجموعه ارزیابی انتخاب شد. شرایط بهینه ساختار شبکه عصبی بهینه برای هر سه مجموعه داده، در جدول ۲-۱۷ خلاصه شده است. مقادیر فعالیت های دارویی داده های



مجموعه آزمون با استفاده از ساختار شبکه عصبی بهینه مربوط به هر مجموعه داده پیش‌بینی شد. بنابراین مدل شبکه عصبی با نماد LAD-LASSO-LM-ANN (با معماری‌های متفاوت مندرج در جدول ۲-۱۷ برای هر مجموعه داده) به‌عنوان مدل برتر برای پیش‌بینی مقادیر فعالیت دارویی ترکیبات مورد مطالعه در هر سه مجموعه داده انتخاب شد.

جدول ۲-۱۷ ساختارهای شبکه‌های توسعه یافته با توصیف‌کننده‌های منتخب LAD\_LASSO با کمترین MSE مجموعه ارزیابی

مجموعه داده‌ها	تعداد توصیف‌کننده	تابع آموزش	تابع انتقال	تعداد گره	تعداد دور آموزش	MSE <sub>validation</sub>	R <sup>2</sup> <sub>validation</sub>
بازدارنده‌های ایدز	۸	تنظیم بایزین	لگاریتم-سیگموئید	۶	۵	۰/۱۶	۰/۸۳
	۵	لونیبرگ-مارکوارت	لگاریتم-سیگموئید	۶	۵	۰/۱۲	۰/۹۱
	۸	تنظیم بایزین	تانژانت-سیگموئید	۲	۵	۰/۱۶	۰/۸۳
	۵	لونیبرگ-مارکوارت	تانژانت-سیگموئید	۷	۲۰	۰/۱۴	۰/۸۶
بازدارنده‌های سرطان کولورکتال	۱۴	تنظیم بایزین	لگاریتم-سیگموئید	۸	۵۰	۰/۱۹	۰/۸۱
	۱۴	لونیبرگ-مارکوارت	لگاریتم-سیگموئید	۲	۵	۰/۰۹	۰/۹۰
	۱۰	تنظیم بایزین	تانژانت-سیگموئید	۲	۵۰	۰/۱۹	۰/۸۲
	۱۳	لونیبرگ-مارکوارت	تانژانت-سیگموئید	۲	۵	۰/۱۰	۰/۹۱
بازدارنده‌های سرطان ریه	۹	تنظیم بایزین	لگاریتم-سیگموئید	۲	۱۰	۰/۱۰	۰/۸۵
	۷	لونیبرگ-مارکوارت	لگاریتم-سیگموئید	۳	۱۰	۰/۰۹	۰/۸۵
	۷	تنظیم بایزین	تانژانت-سیگموئید	۴	۴۰	۰/۰۹	۰/۸۶
	۷	لونیبرگ-مارکوارت	تانژانت-سیگموئید	۴	۴۰	۰/۰۸	۰/۸۶

به منظور تأیید کارایی روش پیشنهادی LAD-LASSO-LM-ANN در این مطالعه، روش انتخاب متغیر کلاسیک SR نیز بر روی داده‌های مجموعه آموزش و ارزیابی هر سه مجموعه داده‌ها اجرا شد. تعداد توصیف‌کننده‌های منتخب SR برای هر سه مجموعه داده‌ها به ترتیب برابر با تعداد ۱۲، ۱۴ و ۱۱ به دست آمد. از این توصیف‌کننده‌ها برای طراحی و بهینه‌سازی مدل‌های شبکه عصبی استفاده شد و بهینه‌سازی پارامترهای ANN به طور هم‌زمان انجام شد. نتایج نشان داد مدل ANN با توابع آموزش LM برای هر سه مجموعه داده و با استفاده از ۱۲ توصیف‌کننده، ۳ گره در لایه پنهان و ۵ دور آموزشی برای مجموعه داده‌های ضد ایدز، ۱۴ توصیف‌کننده، ۲ گره در لایه پنهان و ۱۰ دور آموزشی برای مجموعه داده‌های ضد سرطان کارسینوم کولورکتال و ۱۱ توصیف‌کننده، ۴ گره در لایه پنهان و ۲۵ دور آموزشی برای مجموعه داده‌های ضد سرطان ریه، کمترین MSE را برای مجموعه ارزیابی ایجاد نمود. مدل‌های بهینه SR-LM-ANN برای پیش‌بینی فعالیت دارویی مجموعه آزمون استفاده شد.

## ۲-۴-۷ ارزیابی مدل LAD-LASSO-LM-ANN هر سه مجموعه داده

هدف اصلی از ساخت مدل QSAR، ایجاد یک مدل با قدرت پیش‌بینی مناسب، قابل اعتماد و دقیق برای پیش‌بینی فعالیت دارویی ترکیبات جدید است. در این تحقیق، قدرت پیش‌بینی، اعتبار و تعمیم‌پذیری مدل توسعه یافته LM-ANN با استفاده از توصیف‌کننده‌های منتخب روش انتخاب متغیر LAD-LASSO با به‌کارگیری پیش‌بینی داده‌های مجموعه آزمون و پیش‌بینی پاسخ کل ترکیبات با استفاده از تکنیک LOO، آزمون‌های پراکندگی هم‌چون Y-تصادفی، دامنه کاربرد و محاسبه پارامترهای آماری مورد ارزیابی قرار گرفت.

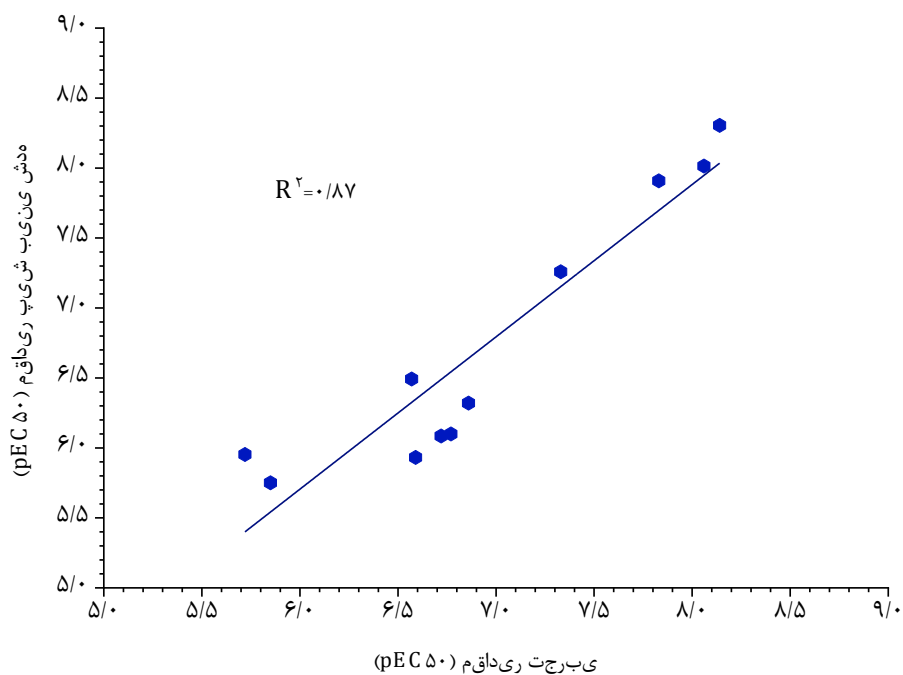
۲-۴-۷-۱ ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از پیش بینی داده‌های مجموعه

## آزمون

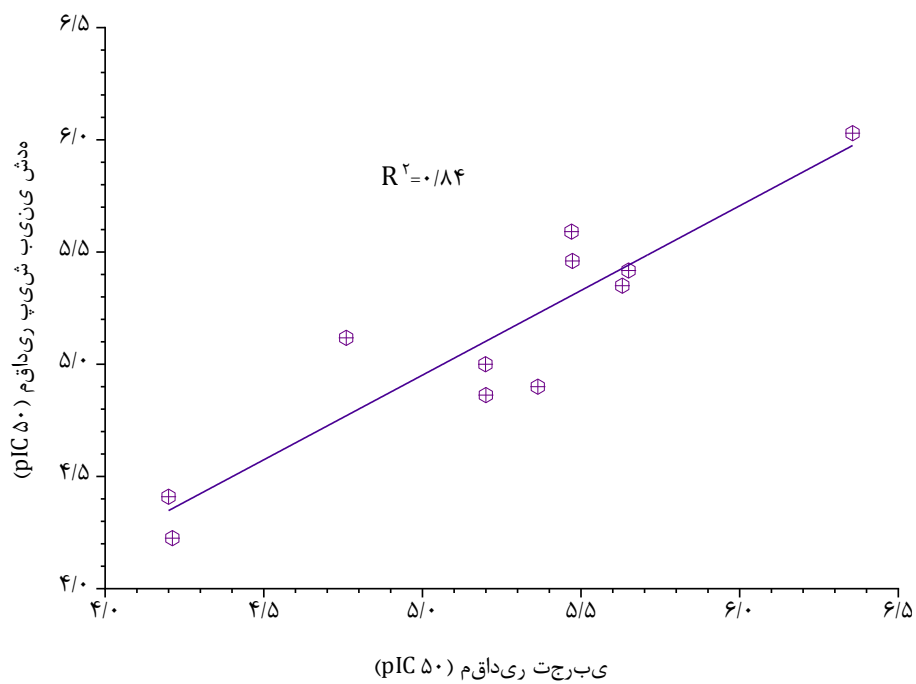
به منظور بررسی اعتبار و قدرت پیش‌بینی مدل بهینه LAD-LASSO-LM-ANN توسعه یافته برای هر مجموعه داده، فعالیت دارویی داده‌های مجموعه آزمون با استفاده از مدل آموزش دیده در شرایط بهینه پیش‌بینی شد. لازم به ذکر است که مجموعه داده‌های آزمون در هیچ یک از مراحل انتخاب توصیف کننده‌های مؤثر و مدل‌سازی حضور نداشته‌اند. مقادیر فعالیت پیش‌بینی شده به وسیله مدل LAD-LASSO-LM-ANN برای ترکیبات موجود در مجموعه آزمون هر سه مجموعه داده در جدول ۲-۱۸ آورده شده است. نتایج حاصل حاکی از خطای کم و قابل قبول اغلب ترکیبات مجموعه آزمون است. برای بررسی بیشتر قدرت پیش‌بینی مدل‌های توسعه یافته LAD-LASSO-LM-ANN مقادیر فعالیت پیش‌بینی شده به وسیله مدل LAD-LASSO-LM-ANN بر حسب فعالیت دارویی تجربی ترکیبات موجود در مجموعه آزمون رسم شد. مقادیر  $R^2$  بزرگ‌تر از  $0/6$  نشان‌دهنده قدرت پیش‌بینی قابل قبول مدل LAD-LASSO-LM-ANN در هر سه مجموعه داده مورد مطالعه است.

جدول ۱۸-۲ نتایج حاصل از ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از مجموعه آزمون

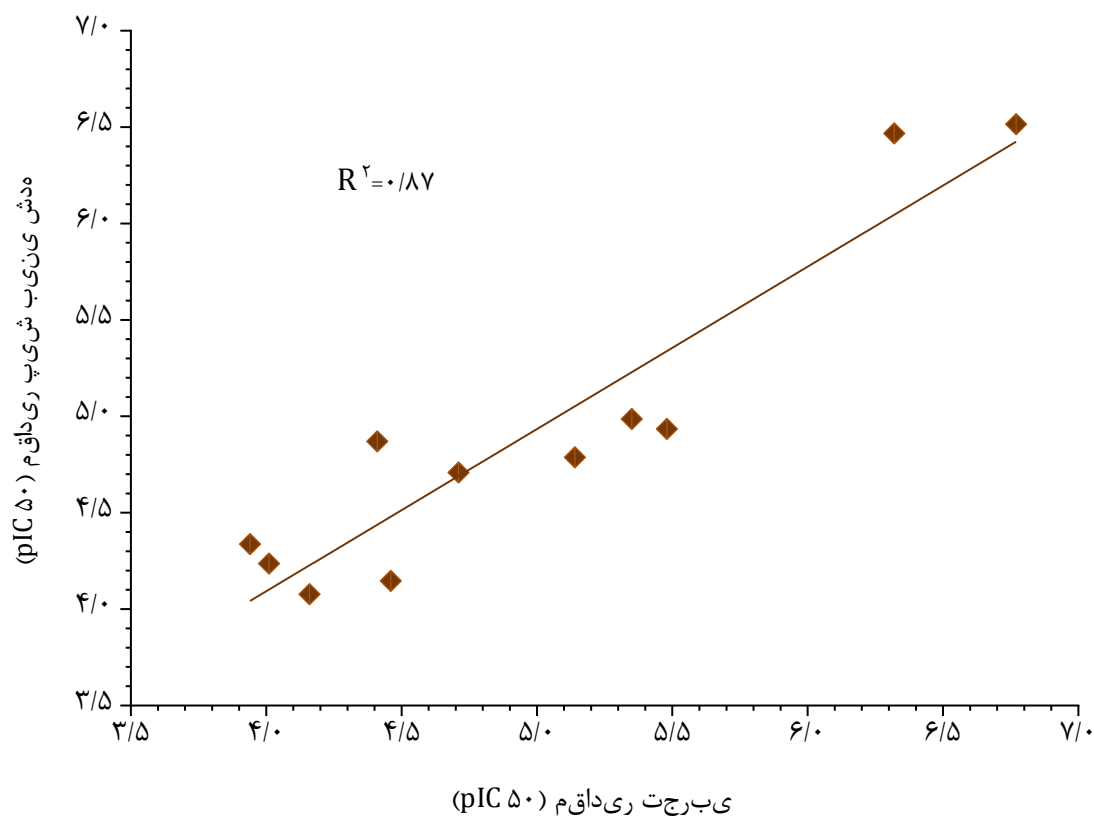
مجموعه داده‌ها	شماره ترکیب	فعالیت دارویی ( $pIC_{50}/pEC_{50}$ )		درصد خطا
		مقدار واقعی	مقدار پیش‌بینی شده	
بازدارنده‌های آیدز	۳	۶/۷۲	۶/۰۹	-۹/۴۴
	۸	۷/۳۳	۷/۲۶	-۰/۹۷
	۱۸	۸/۱۴	۸/۳	۲/۰۳
	۲۴	۷/۸۳	۷/۹۱	۱/۰۰
	۲۵	۸/۰۶	۸/۰۱	-۰/۵۷
	۳۱	۶/۵۷	۶/۴۹	-۱/۱۸
	۳۵	۶/۷۷	۶/۱	-۹/۹
	۴۴	۶/۸۶	۶/۳۲	-۷/۸۷
	۵۸	۶/۵۹	۵/۹۳	-۹/۹۸
	۶۳	۵/۷۲	۵/۹۵	۴/۰۷
بازدارنده‌های سرطان کارسینوم کولورکتال	۲	۶/۳۶	۶/۰۳	-۵/۱۴
	۵	۵/۴۷	۵/۴۶	-۰/۲۴
	۹	۵/۳۶	۴/۹	-۸/۶۵
	۱۰	۵/۲	۵	-۳/۸۳
	۱۱	۵/۶۳	۵/۳۵	-۴/۹۷
	۱۶	۴/۲۱	۴/۲۳	-۰/۳۱
	۱۹	۴/۲	۴/۴۱	۴/۹۹
	۲۹	۵/۴۷	۵/۵۹	۲/۲۱
	۳۸	۴/۷۶	۵/۱۲	۷/۵۱
	۴۸	۵/۶۵	۵/۴۲	-۴/۱۱
بازدارنده‌های سرطان ریه	۱	۵/۱۴	۴/۷۹	-۶/۸۷
	۷	۵/۴۸	۴/۹۳	-۹/۹۶
	۱۹	۳/۹۴	۴/۳۴	۱۰/۰۹
	۲۲	۴/۷۱	۴/۷۱	-۰/۰۴
	۲۴	۴/۴۶	۴/۱۵	-۷/۰۴
	۲۶	۴/۰۱	۴/۲۴	۵/۶۵
	۲۸	۴/۱۶	۴/۰۸	-۱/۹۸
	۴۲	۴/۴۱	۴/۸۷	۱۰/۴۳
	۶۲	۵/۳۵	۴/۹۹	-۶/۱۸
	۶۶	۶/۳۲	۶/۴۷	۲/۳۳
۶۷	۶/۷۷	۶/۵۲	-۳/۷۶	



شکل ۲۲-۲ نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی برای داده‌های ضد ایدز مجموعه آزمون



شکل ۲۳-۲ نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی برای داده‌های ضد سرطان کارسینوم کولورکتال مجموعه آزمون



شکل ۲-۲۴ نمودار تغییرات مقادیر پیش‌بینی شده در مقابل مقادیر تجربی برای داده‌های ضد سرطان ریه مجموعه آزمون

همانطور که در بخش ۲-۴-۷ اشاره شد، از روش SR-LM-ANN برای بررسی عملکرد مدل LAD-LASSO-LM-ANN استفاده شد. پاسخ ترکیبات مجموعه آزمون با استفاده از مدل‌های شبکه عصبی بهینه با معماری‌های اشاره شده در بخش ۲-۴-۷ پیش‌بینی شد. مقادیر MSE مجموعه داده آزمون به ترتیب برابر با ۰/۳۲، ۰/۲۶ و ۰/۴۶ برای مجموعه داده‌های ایدز، سرطان کارسینوم کولورکتال و سرطان ریه به دست آمد. مقادیر  $R^2$  مربوط به مجموعه آزمون نیز به ترتیب برابر با ۰/۷۶، ۰/۷۴ و ۰/۷۶ برای مجموعه داده‌های ایدز، سرطان کارسینوم کولورکتال و سرطان ریه به دست آمد. نتایج نشان می‌دهد که مدل LAD-LASSO-LM-ANN در پیش‌بینی فعالیت‌های دارویی ترکیبات مورد نظر با قدرت بیش‌تری نسبت به مدل SR-LM-ANN عمل کرده است.

۲-۴-۷-۲ ارزیابی مدل LAD-LASSO-LM-ANN با پیش بینی فعالیت دارویی تمام ترکیبات

### مجموعه داده‌ها با استفاده از روش رد مرحله‌ای تک تک

در این بخش برای ارزیابی مدل بهینه و بررسی قدرت پیش بینی مدل از تکنیک قدرتمند دیگری به‌عنوان روش رد مرحله‌ای تک تک (LOO) استفاده شد. با به‌کارگیری این تکنیک هر داده مورد مطالعه یک‌بار به‌عنوان داده آزمون خارج شد و مدل LAD-LASSO-LM-ANN با باقی‌مانده داده‌ها آموزش داده شد و فعالیت دارویی داده خارج شده با استفاده از مدل آموزش دیده، پیش‌بینی شد. این عملیات برای همه ترکیبات موجود در مجموعه داده‌ها تکرار شد و فعالیت دارویی همه داده‌ها توسط مدل LAD-LASSO-LM-ANN پیش‌بینی شدند. نتایج حاصل از پیش‌بینی همه داده‌ها بر اساس تکنیک LOO، برای هر سه مجموعه داده‌ها در جدول ۲-۱۹، جدول ۲-۲۰ و جدول ۲-۲۱ آورده شده‌اند. برای مطالعه بیشتر قدرت پیش بینی مدل، نمودار مقادیر پیش‌بینی شده فعالیت دارویی تمام ترکیبات بر حسب مقادیر تجربی فعالیت دارویی آن‌ها رسم شد. مقادیر  $Q_{LOO}^2$  برای هر سه مجموعه داده‌ها، از مقدار قابل قبول  $0/5$  بزرگ‌تر است، که نشان‌دهنده پایداری و استحکام مناسب مدل‌های بهینه توسعه یافته با توصیف‌کننده‌های منتخب روش LAD-LASSO می‌باشد. علاوه بر این نمودار باقی‌مانده‌ها برای هر سه مجموعه داده‌ها، از رسم مقادیر پیش‌بینی شده فعالیت دارویی بر حسب مقادیر واقعی آن‌ها به‌دست آمد (شکل ۲-۲۶، شکل ۲-۲۸ و شکل ۲-۳۰). توزیع یکنواخت داده‌ها حول محور صفر، نشان‌دهنده عدم وجود خطای سیستماتیک در مدل‌های شبکه عصبی توسعه یافته با استفاده از روش LAD-LASSO به‌عنوان روش انتخاب متغیر (LAD-LASSO-LM-ANN) است.

جدول ۱۹-۲ نتایج حاصل از ارزیابی مدل LAD-LASSO-LM-ANN با تکنیک LOO برای کل مجموعه داده‌های ضد آیدز

شماره ترکیب	pEC <sub>۵۰</sub>			شماره ترکیب	pEC <sub>۵۰</sub>		
	مقدار واقعی	مقدار پیش‌بینی شده	درصد خطا		مقدار واقعی	مقدار پیش‌بینی شده	درصد خطا
۱	۶/۸۵	۶/۶۲	-۳/۴۱	۳۸	۷/۴۷	۶/۸۶	-۸/۱۱
۲	۷/۴۲	۷/۱۳	-۳/۹۱	۳۹	۶/۱۴	۶/۵۱	۶/۰۳
۳	۶/۷۲	۶/۳۲	-۶/۰۲	۴۰	۴/۵۲	۴/۹	۸/۴۱
۴	۶/۸۹	۷/۴	۷/۴	۴۱	۴/۵۲	۴/۹	۸/۴۱
۵	۶/۶۲	۶/۵	-۱/۸۱	۴۲	۵/۴۱	۶/۴۶	۱۹/۳۶
۶	۶/۸۹	۵/۹۱	-۱۴/۱۷	۴۳	۵/۵۲	۵/۹۵	۷/۷۹
۷	۷/۲۴	۷/۰۲	-۳/۱	۴۴	۶/۱۴	۵/۴۹	-۱۰/۵۹
۸	۷/۳۳	۷/۲۴	-۱/۲۳	۴۵	۶/۸۶	۶/۳۱	-۸/۰۶
۹	۷/۱۱	۶/۴۱	-۹/۹	۴۶	۶/۷	۶/۵۱	-۲/۸۸
۱۰	۸/۳۱	۷/۹۹	-۳/۸۴	۴۷	۷/۱	۶/۵۶	-۷/۶۳
۱۱	۷/۹۴	۷/۵۴	-۵/۰۳	۴۸	۶/۸۱	۶/۲۷	-۷/۸۶
۱۲	۷/۶۹	۸/۰۴	۴/۵۷	۴۹	۶/۷۴	۶/۳۷	-۵/۴۷
۱۳	۸/۲۴	۷/۷۵	۶/۰۰	۵۰	۷/۴۷	۶/۹۳	-۷/۲۱
۱۴	۸/۲۶	۷/۸	-۵/۵۹	۵۱	۶/۴	۶/۷۷	۵/۸۱
۱۵	۷/۹۹	۸/۶۷	۸/۵	۵۲	۶/۶۶	۶/۳۶	-۴/۵۳
۱۶	۸/۱	۸/۰۱	-۱/۰۷	۵۳	۶/۸۴	۶/۸۸	۰/۵۲
۱۷	۸/۲۵	۷/۸۷	-۴/۵۵	۵۴	۶/۲۷	۶/۶۹	۶/۷۷
۱۸	۸/۱۴	۸/۰۵	-۱/۱۶	۵۵	۶/۱۱	۶/۲	۱/۴۷
۱۹	۷/۷۸	۷/۸۶	۱/۰۷	۵۶	۵/۸۱	۶/۴۵	۱۱
۲۰	۸/۳۴	۸/۳۱	-۰/۳۳	۵۷	۵/۲۹	۵/۷۱	۷/۸۷
۲۱	۸/۱۶	۸/۰۲	-۱/۷۵	۵۸	۵/۸	۶/۲۱	۷/۰۷
۲۲	۸/۲۲	۸/۰۵	-۲/۰۴	۵۹	۶/۵۹	۶/۲۵	-۵/۲
۲۳	۷/۸۸	۸/۰۴	۱/۹۷	۶۰	۶/۲	۶/۴۲	۳/۵۲
۲۴	۷/۸۳	۷/۹۲	۱/۱	۶۱	۶/۶۸	۶/۵۵	-۱/۹
۲۵	۸/۰۶	۸/۰۲	-۰/۴۵	۶۲	۵/۹۲	۶/۵۶	۱۰/۸۸
۲۶	۷/۹۷	۸/۳۶	۴/۹	۶۳	۴/۸۴	۵/۲	۷/۴۴
۲۷	۷/۹۲	۸/۲۹	۴/۷۳	۶۴	۵/۷۲	۶/۳۳	۱۰/۶۴
۲۸	۶/۲۴	۶/۵۲	۴/۴۱	۶۵	۵/۵۵	۵/۰۵	-۹/۰۳
۲۹	۷/۲۵	۶/۵۲	-۱۰/۰۹	۶۶	۵/۹۴	۶/۲۳	۴/۸
۳۰	۶/۹۶	۶/۴۵	-۷/۳۸	۶۷	۵/۴۷	۶/۲۷	۱۴/۶۵
۳۱	۶/۵۷	۶/۸۸	۴/۷۳	۶۸	۶/۰۸	۶/۳	۳/۶۶
۳۲	۶/۸	۶/۳۵	-۶/۶۱	۶۹	۶/۱۵	۶/۲۴	۱/۴۹
۳۳	۵/۳۷	۵/۷۷	۷/۴۵	۷۰	۵/۸۵	۵/۷۲	-۲/۱۹
۳۴	۶/۲۱	۶/۵۱	۴/۷۹	۷۱	۵/۸۵	۶/۲۶	۶/۹۹
۳۵	۶/۷۷	۶/۹	۱/۹۲	۷۲	۶/۴۳	۶/۳۲	-۱/۷۵
۳۶	۴/۷۷	۵/۹۳	۲۴/۳	۷۳	۶/۸۲	۶/۲۷	-۸/۰۹
۳۷	۶/۹۶	۶/۴۶	-۷/۲۲				



جدول ۲-۲ نتایج حاصل از ارزیابی مدل LAD-LASSO-LM-ANN به روش رد مرحله‌ای تک تک برای کل مجموعه داده‌های

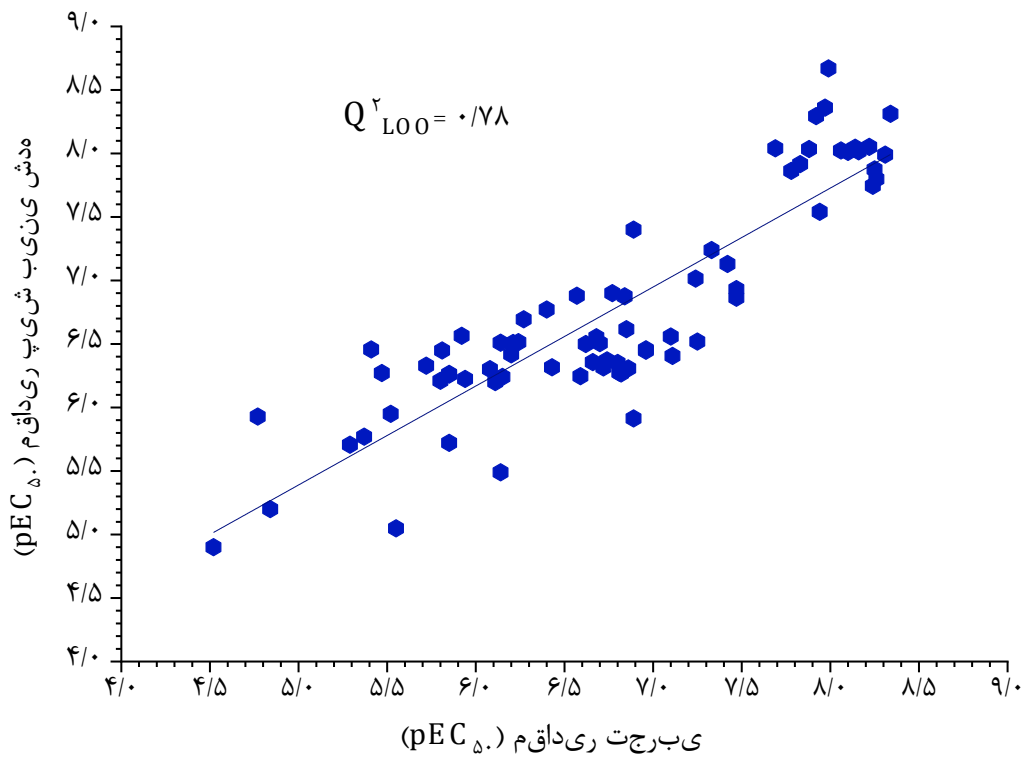
ضد سرطان کارسینوم کولورکتال

شماره ترکیب	pIC <sub>50</sub>		درصد خطا	شماره ترکیب	pIC <sub>50</sub>		درصد خطا
	مقدار واقعی	مقدار پیش‌بینی شده			مقدار واقعی	مقدار پیش‌بینی شده	
۱	۵/۶۷	۵/۴۷	-۳/۵۵	۴۱	۴/۷۴	۴/۸۹	۳/۱۶
۲	۶/۳۶	۵/۹۴	-۶/۵۷	۴۲	۴/۷۲	۴/۷۹	۱/۴۸
۳	۶/۳۲	۵/۹۳	-۶/۱۵	۴۳	۶/۶۷	۶/۱۷	-۷/۴۹
۴	۶/۲۴	۵/۶۹	-۸/۸۴	۴۴	۵/۷۶	۴/۸۴	-۱۶/۰۶
۵	۵/۴۷	۵/۹۶	۸/۹۳	۴۵	۴/۷۹	۴/۷۲	-۱/۵۵
۶	۵/۶۳	۵/۹۱	۴/۹	۴۶	۵/۲۹	۵/۶۹	۷/۵
۷	۵/۸۷	۴/۸۱	۱۸/۰۰	۴۷	۵/۳۹	۴/۸۵	-۱۰/۰۵
۸	۵/۳۶	۵/۸	۸/۱۷	۴۸	۵/۶۵	۵/۴۱	-۴/۲۶
۹	۵/۳۶	۵/۰۳	-۶/۲۳	۴۹	۶/۴۴	۶/۵۵	۱/۷۷
۱۰	۵/۲	۴/۹۴	-۴/۹۳	۵۰	۵/۶۵	۵/۴۸	-۳/۰۵
۱۱	۵/۶۳	۶/۲۷	۱۱/۲۹	۵۱	۷/۰۰	۶/۹۶	-۰/۶۲
۱۲	۷/۰۹	۶/۶۴	-۶/۳۸	۵۲	۶/۲	۶/۴۲	۳/۵۹
۱۳	۶/۷۲	۷/۳۲	۸/۸۹	۵۳	۶/۰۰	۶/۵۲	۸/۶۲
۱۴	۴/۲۳	۴/۴۵	۵/۳۲	۵۴	۴/۳۷	۴/۵	۲/۹۲
۱۵	۴/۲۸	۴/۱۸	-۲/۲۹	۵۵	۵/۰۹	۵/۶۱	۱۰/۲۱
۱۶	۴/۲۱	۴/۱۹	-۰/۵۶	۵۶	۵/۲۰	۵/۱۵	-۱/۰۳
۱۷	۴/۱۸	۴/۴	۵/۱۷	۵۷	۵/۱۷	۶/۰۰	۱۶/۰۵
۱۸	۴/۳۷	۴/۲۲	-۳/۳۳	۵۸	۴/۴۸	۴/۷۲	۵/۲۸
۱۹	۴/۲	۴/۳۶	۳/۹۲	۵۹	۴/۵۶	۵/۲۶	۱۵/۳۵
۲۰	۴/۴۵	۴/۲۴	-۴/۶۱	۶۰	۴/۷۱	۴/۹۲	۴/۵۶
۲۱	۴/۵۴	۵/۰۱	۱۰/۲۶	۶۱	۴/۹۳	۵/۳۹	۹/۲۷
۲۲	۴/۶۵	۴/۷۳	۱/۶۳	۶۲	۵/۳۸	۶/۱۷	۱۴/۷۲
۲۳	۴/۳	۴/۳۴	-۰/۹۵	۶۳	۵/۹۸	۴/۷۴	-۲۰/۷۹
۲۴	۴/۳۵	۴/۳۵	-۰/۰۳	۶۴	۵/۵۷	۴/۹۲	-۱۱/۷
۲۵	۴/۵۴	۴/۶	۱/۳	۶۵	۶/۴۴	۵/۲۴	-۱۸/۷۱
۲۶	۴/۲	۴/۳۵	۳/۶۹	۶۶	۶/۹۲	۷/۶۷	۱۰/۸۵
۲۷	۴/۱۷	۴/۳	۳/۰۹	۶۷	۷/۰۹	۵/۷۴	-۱۹/۰۹
۲۸	۴/۲۱	۴/۲۳	-۰/۳۷	۶۸	۷/۱۵	۶/۴۴	-۹/۹۷
۲۹	۵/۴۷	۴/۹۸	-۹/۰۲	۶۹	۴/۰۸	۴/۱۵	۱/۶۵
۳۰	۴/۳	۴/۴۳	۲/۹۱	۷۰	۴/۱۷	۴/۱۲	-۱/۱۲
۳۱	۴/۲۸	۴/۸۲	۱۲/۶۲	۷۱	۴/۰۱	۴/۴۶	۱۱/۰۸
۳۲	۴/۷۴	۴/۹	۳/۳۴	۷۲	۴/۰۷	۴/۳۳	۶/۴۳
۳۳	۴/۵۹	۴/۱۹	-۸/۷۵				
۳۴	۴/۷۱	۴/۴۷	-۵/۰۹				
۳۵	۵/۱	۵/۱۸	۱/۵۴				
۳۶	۴/۴۹	۵/۲	۱۵/۸۴				
۳۷	۵/۳۴	۴/۵۱	-۱۵/۴۹				
۳۸	۴/۷۶	۴/۶۴	-۲/۴۴				
۳۹	۴/۷۴	۴/۹۴	۴/۲۲				
۴۰	۴/۶۶	۴/۷۵	۲/۰۳				

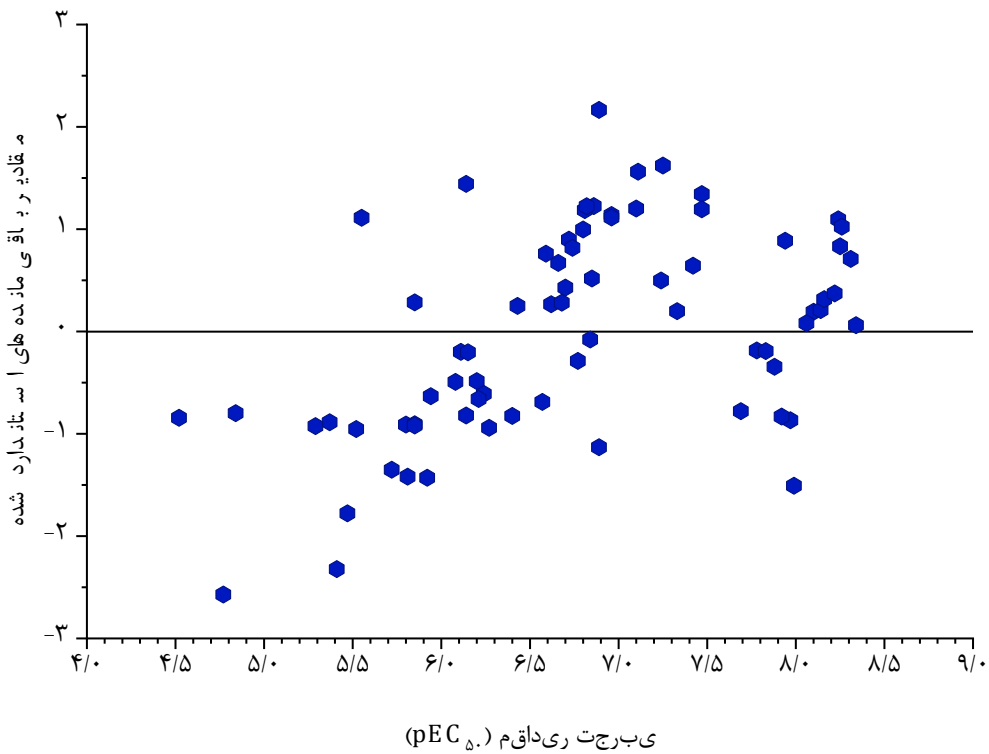
جدول ۲۱-۲ نتایج حاصل از ارزیابی مدل LAD-LASSO-LM-ANN به روش رد مرحله‌ای تک تک برای کل مجموعه داده‌های

ضد سرطان ریه

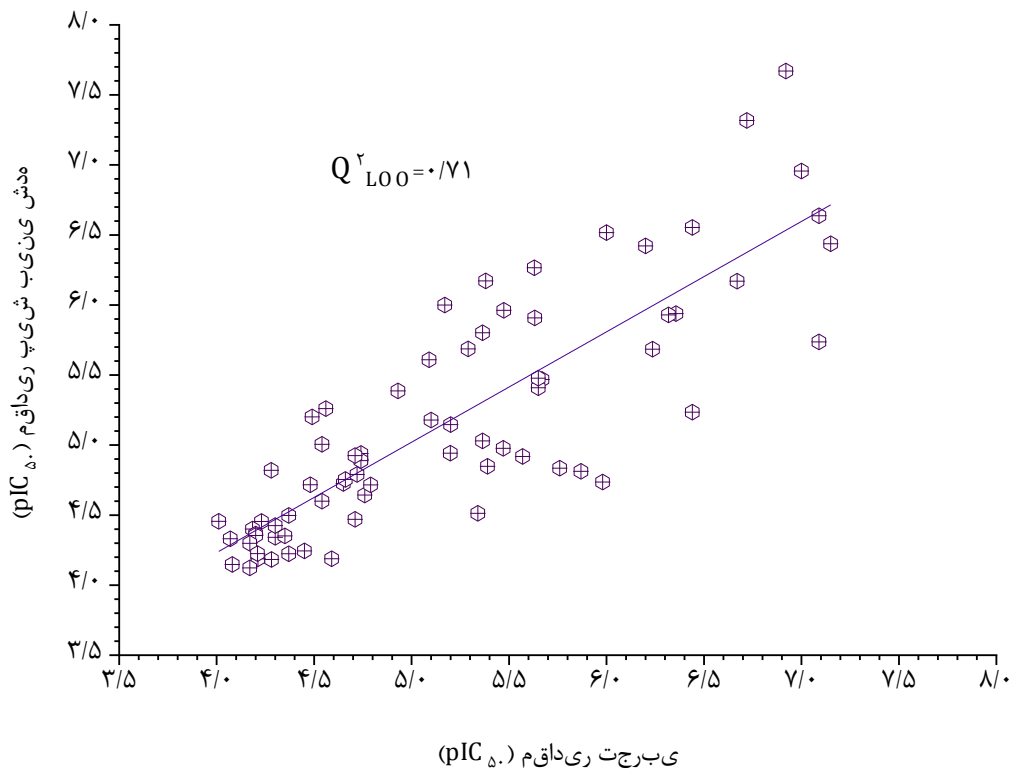
شماره ترکیب	pIC <sub>50</sub>		درصد خطا	شماره ترکیب	pIC <sub>50</sub>		درصد خطا
	مقدار واقعی	مقدار پیش‌بینی شده			مقدار واقعی	مقدار پیش‌بینی شده	
۱	۵/۱۴	۵/۰۸	-۱/۱۲	۴۱	۴/۵۵	۴/۶۲	۱/۴۷
۲	۶/۰۷	۶/۰۷	-۰/۰۲	۴۲	۴/۴۱	۴/۳	-۲/۴۹
۳	۶/۱۱	۶/۳۹	۴/۵۶	۴۳	۶/۲۱	۶/۲۴	۰/۴۸
۴	۵/۷۵	۴/۸۳	-۱۵/۹۸	۴۴	۵/۶۱	۵/۵۹	-۰/۳۶
۵	۴/۹۹	۴/۲۴	-۱۵/۰۶	۴۵	۴/۴۳	۵/۲۳	۱۸/۰۶
۶	۵/۶۵	۵/۷۷	۲/۱۲	۴۶	۵/۰۷	۴/۴۵	-۱۲/۲۱
۷	۵/۴۸	۵/۰۷	-۷/۵۲	۴۷	۴/۹۸	۴/۷۴	-۴/۷۳
۸	۵/۰۳	۵/۵۲	۹/۷۷	۴۸	۵/۲	۵/۳۱	۲/۱۴
۹	۵/۲	۴/۷۷	-۸/۳۴	۴۹	۶/۱۱	۵/۶۸	-۷/۰۴
۱۰	۴/۹۱	۴/۲۸	-۱۲/۸۳	۵۰	۵/۲	۵/۴۶	۴/۹۳
۱۱	۵/۳۷	۵/۹	۹/۹۳	۵۱	۶/۵۱	۶/۷۹	۴/۳
۱۲	۶/۵۵	۶/۸۹	۵/۱۹	۵۲	۵/۶۵	۵/۴۳	-۳/۸۹
۱۳	۶/۵۱	۶/۵۵	۰/۶۱	۵۳	۵/۴۹	۵/۳۱	-۳/۲۸
۱۴	۴/۰۴	۴/۱۱	۱/۶۱	۵۴	۴/۲۲	۴/۵۴	۷/۵۲
۱۵	۴/۰۷	۴/۱	۰/۷۷	۵۵	۵/۰۱	۴/۳۳	-۱۳/۶۳
۱۶	۴/۰۵	۴/۱۳	۱/۹۱	۵۶	۵/۰۳	۴/۴۴	-۱۱/۸۳
۱۷	۴/۱۷	۴/۳۸	۴/۹۸	۵۷	۵/۱۱	۴/۷	-۸/۰۶
۱۸	۴/۱۴	۴/۱۲	-۰/۵۵	۵۸	۴/۳	۴/۸	۱۱/۶۳
۱۹	۳/۹۴	۴/۴۲	۱۲/۱۹	۵۹	۴/۴۹	۴/۸۴	۷/۸
۲۰	۴/۲۳	۴/۱۳	-۲/۴۳	۶۰	۴/۶۶	۴/۶۲	-۰/۸۶
۲۱	۴/۶۱	۴/۷۱	۲/۱۹	۶۱	۴/۸۶	۴/۸۷	۰/۱۶
۲۲	۴/۷۱	۴/۷۴	۰/۷۱	۶۲	۵/۳۵	۴/۷۴	-۱۱/۳۵
۲۳	۴/۱۱	۴/۲۸	۴/۰۵	۶۳	۶/۰۱	۵/۴۸	-۸/۷۹
۲۴	۴/۴۶	۴/۲۵	-۴/۶۱	۶۴	۵/۷۲	۴/۸۶	-۱۵/۰۹
۲۵	۴/۵	۴/۷۴	۵/۳۱	۶۵	۶/۴۶	۵/۸۴	-۹/۵۹
۲۶	۴/۰۱	۴/۳۷	۸/۹۷	۶۶	۶/۳۲	۶/۳۲	۰/۰۰
۲۷	۴/۱	۴/۲۱	۲/۷	۶۷	۶/۷۷	۶/۶۵	-۱/۷۸
۲۸	۴/۱۶	۴/۱	-۱/۵۶	۶۸	۶/۹۲	۶/۰۱	-۱۳/۱۳
۲۹	۵/۱۸	۵/۳۳	۲/۹۵	۶۹	۴/۰۵	۴/۱۸	۳/۱۸
۳۰	۴/۰۷	۴/۲۳	۴/۰۱	۷۰	۴/۱	۴/۲۱	۲/۵۶
۳۱	۴/۲۱	۴/۶۷	۱۰/۹۹				
۳۲	۴/۷۵	۴/۶۵	-۲/۲				
۳۳	۴/۴۹	۴/۳۸	-۲/۳۹				
۳۴	۴/۳۴	۴/۶۷	۷/۵۲				
۳۵	۵/۰۳	۵/۷۱	۱۳/۵۴				
۳۶	۴/۵۶	۴/۶	۰/۸۵				
۳۷	۴/۸۸	۴/۳۵	-۱۰/۸۸				
۳۸	۴/۴۵	۴/۶۹	۵/۳۹				
۳۹	۴/۵	۴/۹۸	۱۰/۶۸				
۴۰	۴/۴۷	۴/۲۹	-۳/۹۹				



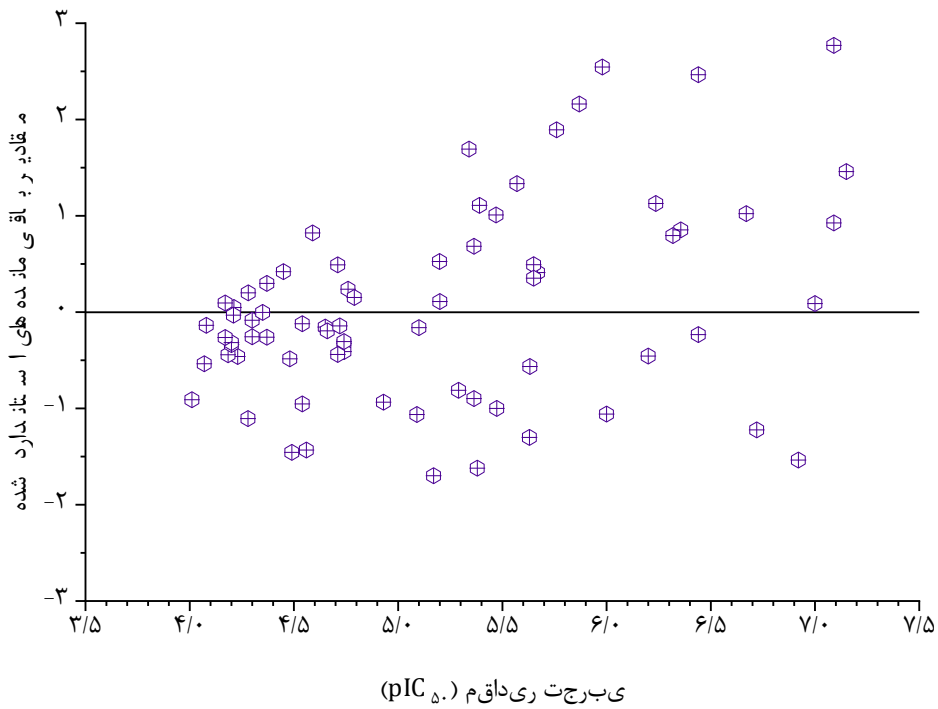
شکل ۲-۲۵ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌های ضد ایدز بر اساس تکنیک LOO در مقابل مقادیر تجربی



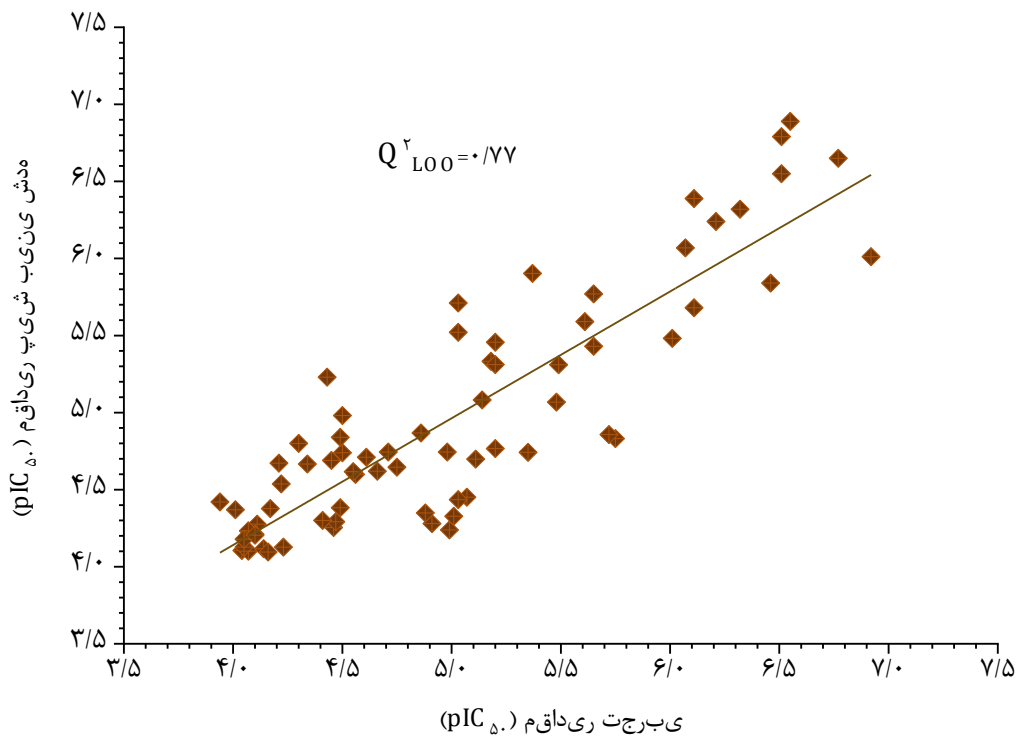
شکل ۲-۲۶ نمودار باقی‌مانده‌های پیش‌بینی شده داده‌های ضد ایدز با استفاده از تکنیک LOO بر حسب مقادیر تجربی



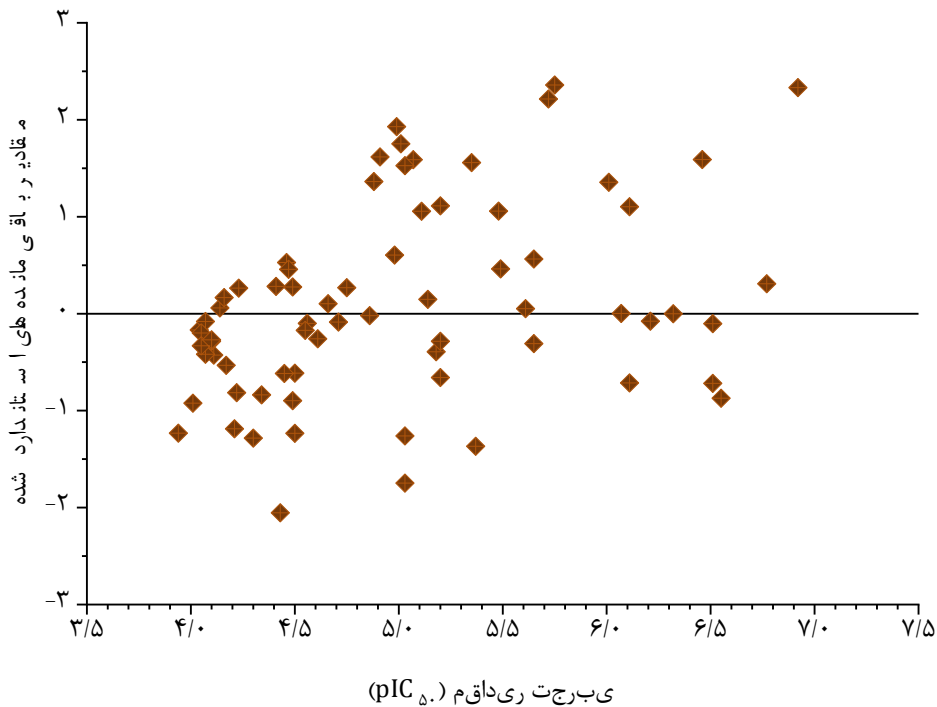
شکل ۲۷-۲ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌های ضد سرطان کارسینوم کولورکتال بر اساس تکنیک LOO در مقابل مقادیر تجربی



شکل ۲۸-۲ نمودار باقی‌مانده‌های پیش‌بینی شده داده‌های ضد سرطان کارسینوم کولورکتال با استفاده از تکنیک LOO برحسب مقادیر تجربی



شکل ۲-۲۹ نمودار تغییرات مقادیر پیش‌بینی شده همه داده‌های ضد سرطان ریه بر اساس تکنیک LOO در مقابل مقادیر تجربی



شکل ۲-۳۰ نمودار باقی‌مانده‌های پیش‌بینی شده داده‌های ضد سرطان ریه با استفاده از تکنیک LOO برحسب مقادیر تجربی

## ۲-۴-۷-۳ ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از پارامترهای آماری

علاوه بر تکنیک‌های ذکر شده، یکی از روش‌های ارزیابی مدل LAD-LASSO-ANN توسعه یافته، محاسبه پارامترهای آماری متفاوت می‌باشد. از جمله پارامترهای آماری پر کاربرد می‌توان به پارامترهای متفاوت بررسی خطای مدل، پارامترهای آماری تروپشا و روی اشاره کرد. بنابراین پارامترهای آماری معرفی شده در بخش ۱-۵-۸-۴ برای فعالیت دارویی پیش‌بینی شده ترکیبات مجموعه آزمون و فعالیت‌های دارویی پیش‌بینی شده برای کل ترکیبات به روش رد مرحله‌ای تک تک محاسبه و در جدول ۲-۲۲ خلاصه شدند. نتایج حاصل از محاسبات آماری جدول ۲-۲۲ نشان می‌دهد که پارامترهای آماری در محدوده قابل قبول قرار دارند. بنابراین مطابق با بخش ۱-۵-۸-۴، با توجه به بزرگتر بودن مقادیر پارامترهای تروپشا و روی همچون  $R_0^2$ ،  $R_0^2$  نسبی،  $R_m^2$  و غیره از مقدار قابل قبول  $0/5$  و نزدیک بودن این پارامترها به مقدار  $R^2$  قدرت پیش‌بینی و تعمیم‌پذیری مدل توسعه یافته LAD-LASSO-ANN اثبات می‌شود. علاوه بر این شیب نمودار حاصل از مقادیر پیش‌بینی شده بر حسب مقادیر تجربی (و بالعکس) در عرض از مبدأ صفر نیز در محدوده  $0/85$  تا  $1/15$  قرار دارند که این نتیجه نیز حاکی از صحت مدل توسعه یافته ANN با توصیف‌کننده‌های منتخب روش LAD-LASSO می‌باشد.

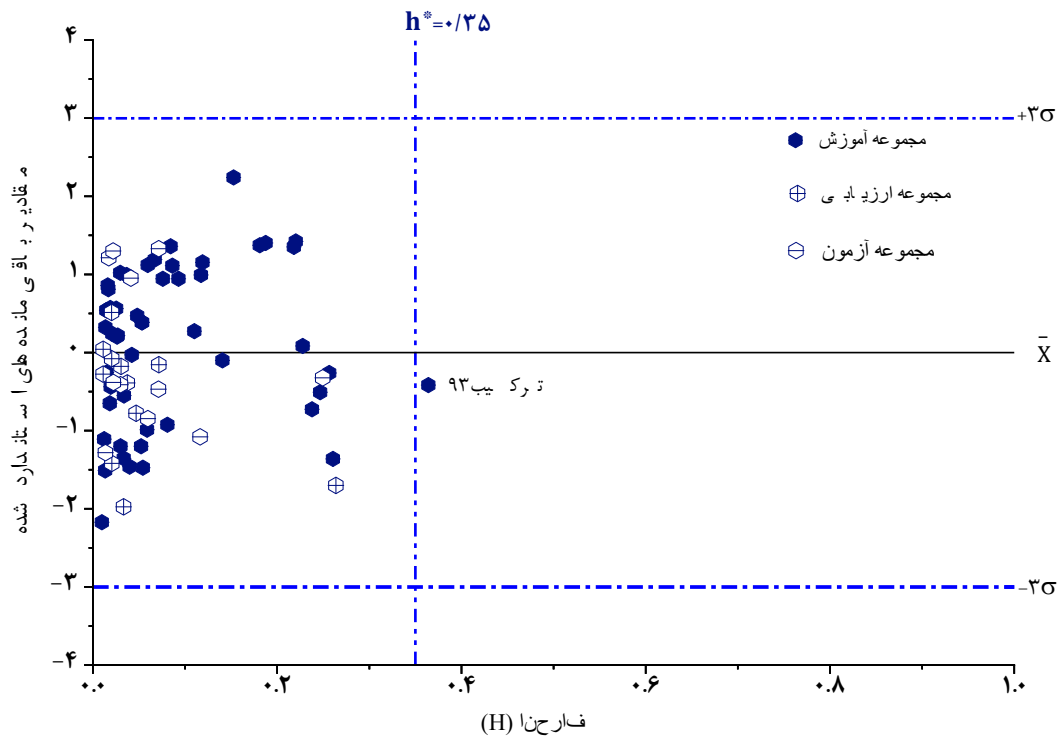
جدول ۲۲-۲ پارامترهای آماری محاسبه شده برای مجموعه آزمون و داده‌های پیش‌بینی شده با تکنیک LOO برای مدل LAD- LASSO-LM-ANN هر سه مجموعه از داده‌ها

ردیف	پارامتر آماری	مجموعه داده‌های ضد سرطان				مجموعه داده‌های ضد سرطان ریه		محدوده قابل قبول
		مجموعه داده‌های ضد ایدز		کارسینوم کولورکتال				
		ترکیبات مجموعه آزمون	کل ترکیبات به روش LOO	ترکیبات مجموعه آزمون	کل ترکیبات به روش LOO	ترکیبات مجموعه آزمون	کل ترکیبات به روش LOO	
۱	PRESS	۱/۴۸	۱۴/۶۲	۰/۸۰	۱۶/۹۶	۱/۱۷	۱۰/۵۷	-
۲	SEP	۰/۳۷	۰/۴۵	۰/۲۷	۰/۴۹	۰/۳۳	۰/۳۹	-
۳	MAE	۰/۲۹	۰/۳۶	۰/۲۳	۰/۳۰	۰/۲۸	۰/۲۸	-
۴	REP(%)	۵/۶۲	۶/۶۱	۵/۱۴	۹/۳۷	۶/۵۴	۷/۷۶	-
۵	MSE	۰/۱۳	۰/۲۰	۰/۰۷	۰/۲۴	۰/۱۱	۰/۱۵	-
۶	MRE	۵/۲۴	۵/۹۲	۴/۴۲	۶/۸۴	۵/۸۹	۵/۹۳	-
۷	R <sup>2</sup>	۰/۸۷	-	۰/۸۴	-	۰/۸۷	-	R <sup>2</sup> > ۰/۶
۸	Q <sub>LOO</sub> <sup>2</sup>	-	۰/۷۸	-	۰/۷۱	-	۰/۷۷	Q <sub>LOO</sub> <sup>2</sup> > ۰/۵
۹	R <sub>0</sub> <sup>2</sup>	۰/۸۶	۰/۷۲	۰/۷۷	۰/۶۶	۰/۸۵	۰/۷۳	نزدیک به R <sup>2</sup>
۱۰	R <sub>0</sub> <sup>2</sup> نسبی	۰/۰۱	۰/۰۸	۰/۰۸	۰/۰۷	۰/۰۲	۰/۰۵	< ۰/۱
۱۱	R <sub>m</sub> <sup>2</sup>	۰/۷۸	۰/۵۹	۰/۶۲	۰/۵۵	۰/۷۵	۰/۶۲	> ۰/۵
۱۲	R <sub>0</sub> <sup>2</sup>	۰/۸	۰/۷۸	۰/۸۴	۰/۷	۰/۸۷	۰/۷۶	نزدیک به R <sup>2</sup>
۱۳	R <sub>0</sub> <sup>2</sup> نسبی	۰/۰۸	۰/۰۰	۰/۰۰	۰/۰۱	۰/۰۰	۰/۰۱	< ۰/۱
۱۴	R <sub>m</sub> <sup>2</sup>	۰/۶۵	۰/۷۲	۰/۵۷	۰/۵۹	۰/۷۳	۰/۶۶	> ۰/۵
۱۵	R-R	۰/۰۶	۰/۰۶	۰/۰۷	۰/۰۱	۰/۰۲	۰/۰۳	< ۰/۳
۱۶	k	۰/۹۷	۱	۰/۹۸	۰/۹۹	۰/۹۸	۱	۰/۸۵ ≤ k ≤ ۱/۱۵
۱۷	k'	۱/۰۳	۱	۱/۰۲	۱	۱/۰۱	۱	۰/۸۵ ≤ k' ≤ ۱/۱۵

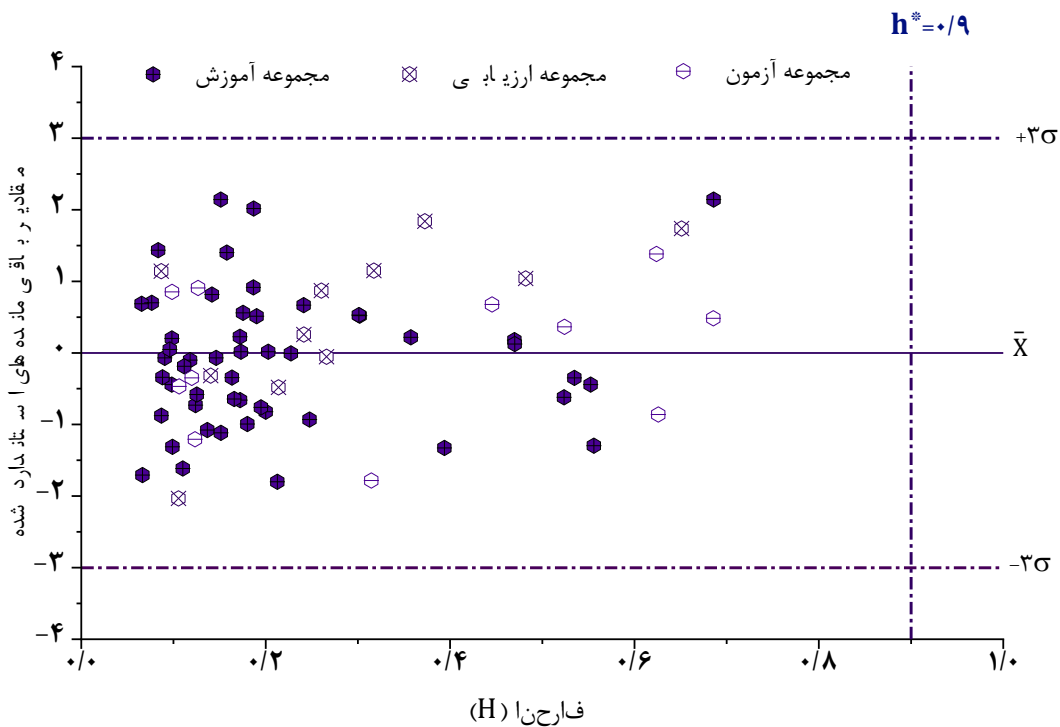
## ۲-۴-۷-۴ ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از دامنه کاربرد

از جمله روش‌های ارزیابی اعتبار مدل QSAR توسعه یافته می‌توان به آزمون دامنه کاربرد اشاره کرد. در واقع با استفاده از این روش، یک فضای شیمیایی تئوری با استفاده از توصیف‌کننده‌های مولکولی مجموعه آموزش و فعالیت دارویی مربوطه ایجاد می‌شود. آزمون دامنه کاربرد مدل توسعه یافته LAD-LASSO-LM-ANN با استفاده از رسم نمودار ویلیام مورد تجزیه و تحلیل قرار گرفت. بنابراین مطابق با توضیحات ارائه شده در بخش‌های ۱-۵-۸-۵ و ۲-۲-۷-۴ مقادیر H هر ترکیب با استفاده از رابطه ۱-۱۳ محاسبه شد. سپس مقادیر باقی‌مانده‌های استاندارد شده طبق رابطه ۱-۱۴ محاسبه شد. نمودار ویلیام از رسم مقادیر باقی‌مانده‌های استاندارد شده بر حسب مقابل مقادیر H به دست آمد. همان‌طور که شکل ۲-۳۱ و شکل ۲-۳۲ نشان می‌دهند، مقادیر محاسبه شده H و باقی‌مانده‌های استاندارد شده برای مجموعه داده‌های ضد ایدز و ضد سرطان کارسینوم کولورکتال، در محدوده قابل قبول قرار دارند و هیچ کدام از داده‌های به کار گرفته شده در ساخت، ارزیابی و آزمون مدل LAD-LASSO-LM-ANN به‌عنوان داده دور افتاده نیستند. با توجه به نمودار ویلیام (شکل ۲-۳۳) مربوط به مجموعه داده‌های ضد سرطان ریه، مشاهده می‌شود که تنها یک داده بزرگ‌تر از مقدار هشدار  $h^*$  است و سایر داده‌ها در محدوده قابل قبول قرار گرفته‌اند. بنابراین مدل‌های LAD-LASSO-LM-ANN توسعه یافته برای هر سه مجموعه داده‌ها از استحکام و اطمینان‌پذیری قابل قبولی برخوردار هستند.

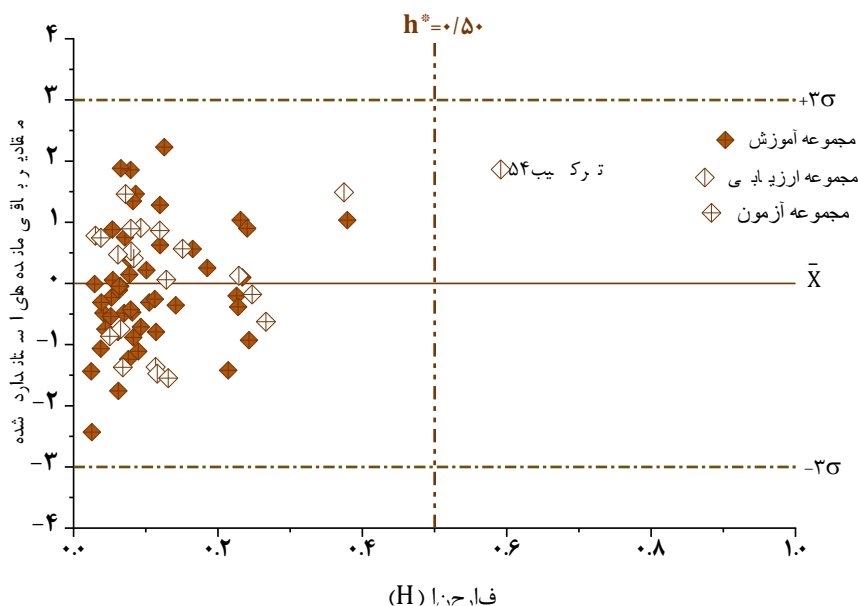




شکل ۳۱-۲ دامنه کاربرد مدل LAD-LASSO-LM-ANN برای مجموعه داده‌های ضد ایدز، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $\pm 3\sigma$  و  $h^*$  است.



شکل ۳۲-۲ دامنه کاربرد مدل LAD-LASSO-LM-ANN برای مجموعه داده‌های ضد سرطان کارسینوم کولورکتال، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $\pm 3\sigma$  و  $h^*$  است.

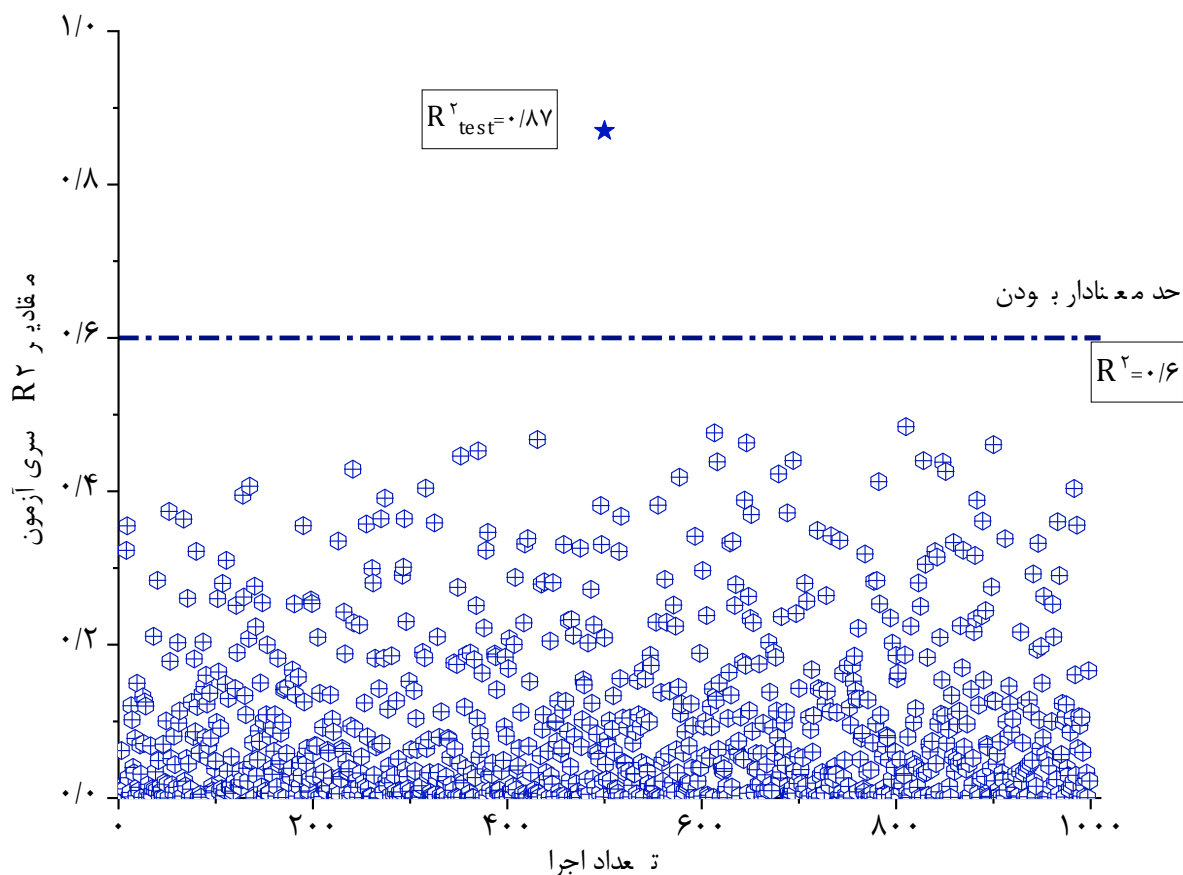


شکل ۲-۳۳ دامنه کاربرد مدل LAD-LASSO-LM-ANN برای مجموعه داده‌های ضد سرطان ریه، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $\pm 3\sigma$  و  $h^*$  است.

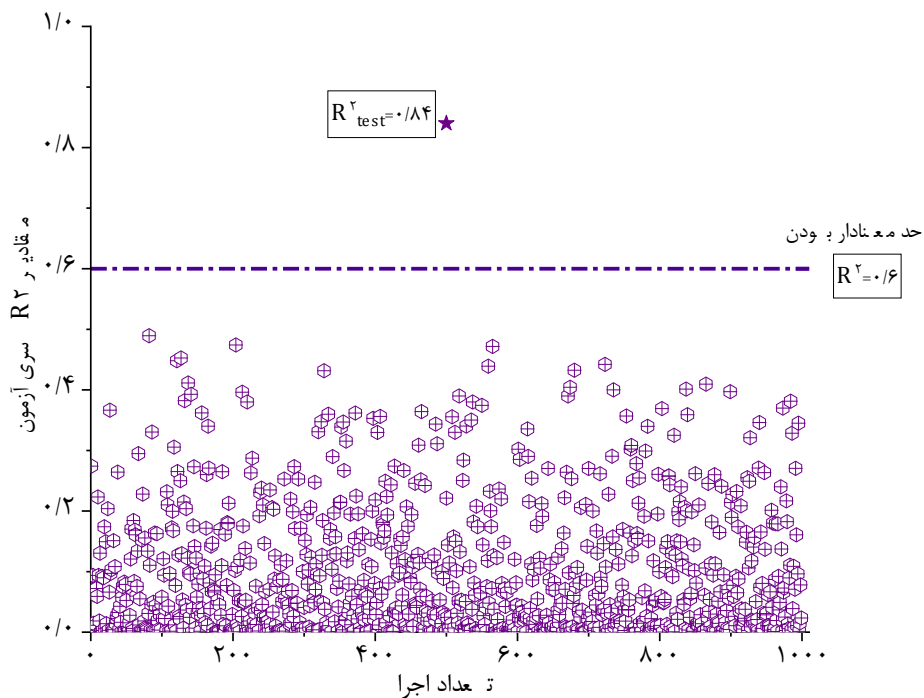
## ۲-۴-۵ ارزیابی مدل LAD-LASSO-LM-ANN با استفاده از آزمون Y-تصادفی

آزمون Y-تصادفی برای بررسی عدم وجود ارتباط تصادفی ایجاد شده توسط مدل LAD-LASSO-LM-ANN بین متغیرهای مستقل و متغیر وابسته به کار گرفته شد. بنابراین برای انجام آزمون Y-تصادفی، ابتدا مقادیر فعالیت دارویی مربوط به هر مجموعه داده، در محدوده تغییرات خودش (حداقل و حداکثر مقدار فعالیت دارویی) ۱۰۰۰ بار به طور کاملاً تصادفی تغییر داده شدند. مدل شبکه عصبی بهینه مربوط به هر مجموعه داده با استفاده از فعالیت‌های دارویی تصادفی آموزش داده شد و مقادیر فعالیت‌های دارویی مجموعه آزمون با استفاده از مدل بهینه LAD-LASSO-LM-ANN توسعه یافته با متغیر وابسته تصادفی، پیش‌بینی شدند. این فرایند برای هر ۳ مجموعه داده انجام شد. نمودار مقادیر پیش‌بینی شده بر حسب مقادیر واقعی برای هر مجموعه داده رسم شد و مقادیر  $R^2$  مربوطه به دست آمد. نتایج  $R^2$  حاصل از پیش‌بینی فعالیت دارویی ترکیبات مجموعه آزمون برای هر ۱۰۰۰ مدل توسعه یافته با متغیر وابسته به دست آمد. مقادیر  $R^2$

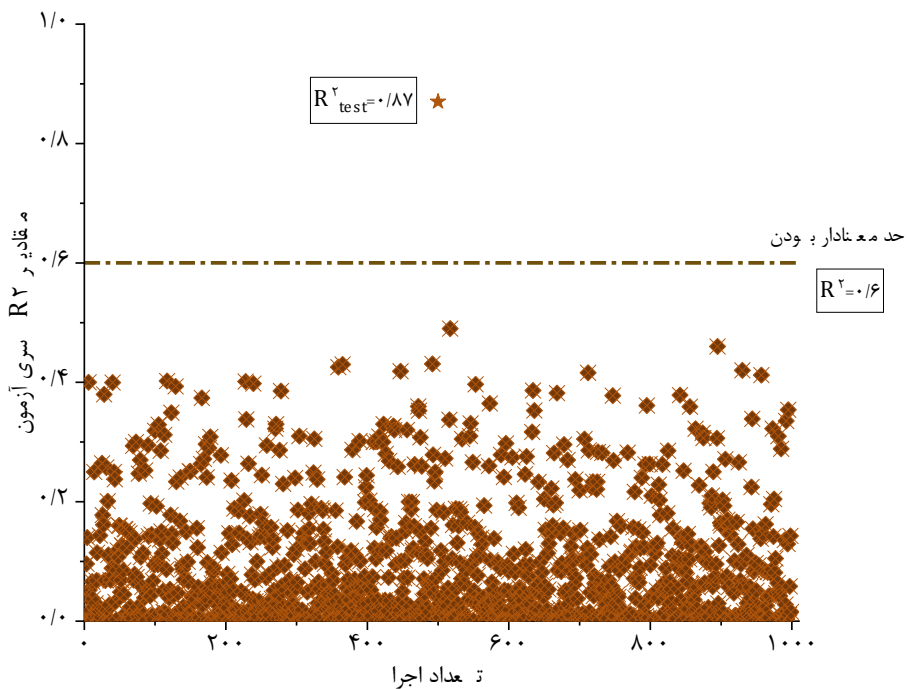
مربوط به ۱۰۰۰ اجرا در هر سه مجموعه داده در شکل ۲-۳۴ تا شکل ۲-۳۵ نشان داده شد. همان‌طور که مشاهده می‌شود مقادیر  $R^2$  حاصل از مدل LAD-LASSO-LM-ANN توسعه یافته با متغیر وابسته تصادفی از مقدار قابل قبول ۰/۶ کوچک‌تر هستند، در نتیجه استنباط می‌شود که مدل‌های توسعه یافته LAD-LASSO-LM-ANN بر اساس ارتباط منطقی بین توصیف‌کننده‌های منتخب روش LAD-LASSO و فعالیت دارویی مربوطه ساخته شده است و ارتباط ایجاد شده به‌طور تصادفی و شانسی به وجود نیامده است.



شکل ۲-۳۴ نمودار مقادیر  $R^2$  به‌دست آمده در آزمون Y-تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای Y-تصادفی و پیش‌بینی فعالیت ترکیبات ضد ایدز مجموعه آزمون به‌وسیله مدل LAD-LASSO-LM-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه



شکل ۳۵-۲ نمودار مقادیر  $R^2$  به دست آمده در آزمون  $Y$ -تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای  $Y$ -تصادفی و پیش‌بینی فعالیت ترکیبات ضد سرطان کارسینوم کولورکتال مجموعه آزمون به وسیله مدل LAD-LASSO-LM-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه



شکل ۳۶-۲ نمودار مقادیر  $R^2$  به دست آمده در آزمون  $Y$ -تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای  $Y$ -تصادفی و پیش‌بینی فعالیت ترکیبات ضد سرطان ریه مجموعه آزمون به وسیله مدل LAD-LASSO-LM-ANN با استفاده از پاسخ تصادفی شده در شرایط بهینه

## ۲-۵ پیش‌بینی شاخص بازداری برخی از ترکیبات آلی فرار با استفاده از

### مدل SCAD-ANN

#### ۲-۵-۱ مقدمه

با توجه به مشکلات و محدودیت‌های کارهای آزمایشگاهی، استفاده از روش‌های تئوری برای محاسبه خواص و ویژگی‌های ترکیبات بسیار حائز اهمیت می‌باشد. امروزه مطالعات کمومتریکس راهی برای استخراج حداکثر اطلاعات از نتایج تجربی با استفاده از انجام یک سری محاسبات آماری و ریاضی می‌باشد. مطالعات ارتباط کمی ساختار - خاصیت (QSPR)، یکی از روش‌های کمومتریکس است که با استفاده از آن می‌توان روابط خطی و یا غیر خطی بین ویژگی‌های ساختاری و خواص ترکیبات را پیدا نمود و از روی آن خاصیت مورد نظر را برای ترکیبات دیگر پیش‌بینی نمود [۱۷۴-۱۸۰]. بنابراین، در شیمی محاسباتی، یافتن رابطه بین خصوصیات ساختاری ترکیبات و خواص آن‌ها مورد توجه محققین است. در میان رویکردهای مختلف شیمی محاسباتی، مطالعه QSPR ابزار مفیدی برای یافتن یک رابطه منطقی مناسب است و به‌طور گسترده در زمینه‌های مختلف مانند کمومتریکس و شیمی دارویی استفاده شده است. QSPR به‌عنوان یک ابزار قوی برای تجزیه و تحلیل خواص کروماتوگرافی نیز شناخته شده است. اخیراً محاسبه شاخص بازداری در مقالات متفاوت گزارش شده است [۸۹، ۹۱، ۱۸۶-۱۸۱]. با توجه به این‌که شناسایی ترکیبات اغلب با تطبیق پیک‌های کروماتوگرافی گازی (GC) آنالیت با پیک‌های استاندارد انجام می‌شود، درحالی‌که نمونه‌های استاندارد خالص گاهی اوقات در دسترس نیستند، امر تطبیق پیک با مشکل مواجه می‌شود. تکنیک GC همچنین به آماده‌سازی نمونه و ستون و بهینه‌سازی پارامتر نیاز دارد که زمان‌بر و پرهزینه است [۹۱، ۱۸۷]. از طرفی، به دلیل پیچیدگی روش‌های تجزیه‌ای و تحلیل کمی، تعیین شاخص بازداری (RI) مواد شیمیایی

---

<sup>1</sup>Retention index

همیشه برای محققین امری زمان بر است. بنابراین، معرفی روش‌های کارآمد برای پیش‌بینی مقادیر RI برای ترکیبات ناشناخته بدون هیچ گونه اندازه‌گیری یا آزمایش مورد توجه است [۸۹]. مدل‌های مطالعه ارتباط کمی ساختار- شاخص بازداری (RI) (QSRR) می‌توانند برای پیش‌بینی شاخص‌های بازداری ساختارهای شیمیایی مورد استفاده قرار گیرند [۸۹, ۹۱, ۱۸۶]. همان‌طور که گفته شد، هدف اصلی مطالعه QSRR، ایجاد یک رابطه ریاضی بین RI به‌عنوان پاسخ یک سیستم کروماتوگرافی و خواص مولکولی به‌عنوان توصیف‌کننده‌های ساختار آنالیت، می‌باشد. بنابراین انتخاب مؤثرترین توصیف‌کننده‌ها با بیش‌ترین ارتباط با متغیر پاسخ، امر مهمی در ساخت مدل QSRR است. از این‌رو، معرفی تکنیک‌های انتخاب متغیر جدید و کارآمد برای افزایش عملکرد مدل‌سازی، بهبود تفسیرپذیری مدل، و کاهش پیچیدگی محاسبات و ذخیره‌سازی ضروری است [۱۸۸, ۱۸۹]. بنابراین، به‌عنوان یک روش انتخاب متغیر قدرتمند، از SCAD برای ساخت مدل QSRR با نتایج رضایت‌بخش استفاده شد [۱۰۷, ۱۰۹, ۱۹۰]. با توجه به جنبه‌های هیجان‌انگیز SCAD، مانند پراکندگی و ثبات بالا، این روش می‌تواند جایگزین مناسبی برای روش‌های کلاسیک انتخاب متغیر و یا LASSO باشد. با توجه به جستجوی انجام شده در مقالات منتشر شده، هیچ گزارشی مبنی بر استفاده از SCAD-ANN در مطالعات QSRR وجود ندارد. بنابراین، در این مطالعه از روش SCAD به‌عنوان تکنیک کاهش متغیر در مطالعه QSRR مقادیر RI برای برخی از ترکیبات آلی فرار استفاده شد.

طبق توضیحات ذکر شده، در این بخش از مطالعه، روش SCAD به‌عنوان یک روش جدید انتخاب متغیر جفت شده با روش ANN به‌عنوان مدل‌سازی غیرخطی قدرتمند (SCAD-ANN) برای پیش‌بینی مقادیر شاخص‌های بازداری (RI) ترکیبات آلی فرار (VOCs) متفاوت استفاده شد. VOC ها ترکیبات شیمیایی آلیفاتیک و معطر با وزن مولکولی و نقطه جوش کم هستند [۱۹۱]. منابع VOC شامل حلال‌ها، سوخت‌ها، رنگ‌ها، مواد شوینده، سیگار و مواد غذایی است. بیش‌تر VOC های تولید شده توسط صنایع

اغلب حاوی بنزن، تولوئن، زایلن، فوران و کلروفرم بوده که برای سلامت انسان مضر هستند [۱۹۲]. در این پژوهش، پس از محاسبه توصیف‌کننده‌های ساختاری برای مجموعه‌های مختلف VOCs و غربالگری متغیرها، توصیف‌کننده‌های معنی‌دار مطابق با  $\lambda_{min}$  استخراج شدند. با توجه به رابطه پیچیده و غیرخطی بین متغیرهای وابسته و مستقل، از ANN برای ایجاد مدل QSRR استفاده شد. عملکرد مدل‌های SCAD-ANN توسعه‌یافته با استفاده از پارامترهای آماری برای پیش‌بینی RI های مجموعه‌های آزمون و کل مجموعه داده‌ها مورد ارزیابی قرار گرفت. با به‌کارگیری SCAD-ANN پیشنهادی مشخص شد که روش مدل‌سازی ANN جفت شده با تکنیک انتخاب متغیر جریمه‌شده SCAD، مدل‌های QSRR دقیقی را برای پیش‌بینی شاخص‌های بازدارندگی ترکیبات VOC ایجاد می‌کند.

## ۲-۵-۲ مجموعه داده‌ها

به‌منظور توسعه مدل‌های QSRR و بررسی کارایی روش انتخاب متغیر SCAD، از دو مجموعه داده متفاوت استفاده شد. اولین مجموعه شامل ۱۳۲ ترکیب آلی فرار متشکل از آلکان‌ها، آلکن‌ها، آمین‌ها، اترها، الکل‌ها، آلکیل بنزن و آلکیل هالیدها بود. شاخص بازدارندگی ترکیبات مورد مطالعه مجموعه داده A با استفاده از دستگاه GC مدل ۴۳۹ Packard Becker (Delft, Netherlands) مجهز به دو آشکارساز هدایت حرارتی اندازه‌گیری شد. فاز ساکن C67 با نام شیمیایی ۱۹،۱۹-Diethyl-14,24-ditridecylheptatricosane بود و داده‌های GC برای حدود ۱۳۲ ترکیب آلی فرار در C67 در دمای  $130^{\circ}\text{C}$  اندازه‌گیری شد [۱۹۳]. مجموعه داده B شامل ۵۲ ترکیب آلی فرار متشکل از پیرازین‌ها، پیریدین‌ها، فوران‌ها و غیره است. RI های ترکیبات مورد مطالعه مجموعه داده دوم با استفاده از دستگاه GC Varian CP-3800 و ستون DB-5ms (30m×0.25mm×0.25mm) استخراج شد [۱۹۴]. ساختار ترکیبات به یک فایل با فرمت سیستم ورودی خطی ورودی مولکولی ساده شده (SMILES) تبدیل شد. نام شیمیایی و فرمت SMILES ساختارهای مورد مطالعه همراه با مقادیر واقعی و پیش‌بینی شده RI مربوط به هر دو مجموعه داده به ترتیب در جدول ۲-۲۳

و جدول ۲-۲۴ خلاصه شده‌اند. مقادیر RI برای همه ترکیبات در مجموعه خود، تحت شرایط یکسان اندازه‌گیری شده است. لازم به ذکر است که دو مجموعه داده، همگنی لازم برای ترکیب شدن به‌عنوان یک مجموعه داده را نداشتند، بنابراین به همین دلیل، به‌عنوان دو مجموعه داده مجزا برای مطالعات بیشتر مورد استفاده قرار گرفتند. مقادیر RI به‌عنوان متغیر وابسته در کل مطالعه در نظر گرفته شد. مجموعه داده‌ها با استفاده از نرم‌افزار R و اجرای الگوریتم KS به مجموعه‌های آموزش، ارزیابی و آزمون تقسیم شدند. در ادامه توضیحات مربوط به ساخت مدل QSRR، مجموعه داده اول شامل ۱۳۲ ترکیب آلی فرار به نام مجموعه داده A و مجموعه داده دوم شامل ۵۲ ترکیب آلی فرار با نام مجموعه داده B نامگذاری شده است.



جدول ۲-۲۳ ساختار شیمیایی و SMILES مربوط به ترکیبات مورد مطالعه مجموعه داده‌های A به همراه شاخص بازداری

ردیف	ساختار شیمیایی	SMILES	RI واقعی	RI پیش بینی شده
۱ <sup>t</sup>	1-Butanol	C(CCC)O	۰/۴۲۲۱	۰/۴۱۴۳
۲	2-Methyl-2-propanol	CC(C)(O)C	۰/۳۳۴۱	۰/۳۳۲۳
۳ <sup>v</sup>	1-Pentanol	C(CCCC)O	۰/۴۹۷۲	۰/۴۹۰۱
۴ <sup>t</sup>	2-Methyl-2-butanol	CC(CC)(O)C	۰/۴۲۵۹	۰/۴۱۸۳
۵	1-Hexanol	C(CCCCC)O	۰/۵۷۰۳	۰/۵۶۷۹
۶ <sup>t</sup>	Cyclohexanol	C1(CCCCC1)O	۰/۶۰۶۲	۰/۶۲۴۱
۷	2-Methyl-2-pentanol	CC(CCC)(O)C	۰/۴۹۲۳	۰/۴۹۳۳
۸	1-Heptanol	C(CCCCCCC)O	۰/۶۴۳	۰/۶۵۱۸
۹ <sup>v</sup>	2-Methyl-2-hexanol	CC(CCCC)(O)C	۰/۵۶۱۱	۰/۵۶۹۵
۱۰	2-Butanol	C[C@@H](CC)O	۰/۳۸۸۵	۰/۳۸۵۵
۱۱	2-Pentanol	C[C@@H](CCC)O	۰/۴۵۹۴	۰/۴۶۳۹
۱۲	2-Hexanol	C[C@@H](CCCC)O	۰/۵۳۱۳	۰/۵۳۵۲
۱۳	2-Heptanol	C[C@@H](CCCCC)O	۰/۶۰۲۷	۰/۶۰۳۷
۱۴ <sup>t</sup>	2-Phenylethanol	C(Cc1ccccc1)O	۰/۷۵۰۴	۰/۷۴۰۴
۱۵ <sup>v</sup>	Benzyl alcohol	C(c1ccccc1)O	۰/۶۹۲۸	۰/۶۵۵۵
۱۶	Pentanal	C(=O)CCCC	۰/۴۶۰۵	۰/۴۶۸۷
۱۷	Hexanal	C(=O)CCCCC	۰/۵۳۳۶	۰/۵۳۳۵
۱۸ <sup>v</sup>	2-Butanone	CC(=O)CC	۰/۳۸۰۸	۰/۳۷۵۴
۱۹ <sup>v</sup>	2-Pentanone	CC(=O)CCC	۰/۴۴۸۱	۰/۴۴۹۴
۲۰	Cyclopentanone	C1(=O)CCCC1	۰/۵۲۷	۰/۵۰۶۵
۲۱	2-Hexanone	CC(=O)CCCC	۰/۵۲۰۶	۰/۵۰۳۵
۲۲ <sup>t</sup>	Cyclohexanone	C1(=O)CCCCC1	۰/۶۰۸	۰/۵۸۵
۲۳	2-Heptanone	CC(=O)CCCCC	۰/۵۹۲۶	۰/۵۶۳۳
۲۴	Dipropylether	O(CCC)CCC	۰/۴۶۸۱	۰/۴۵۴۱
۲۵ <sup>v</sup>	Dibutyl ether	O(CCCC)CCCC	۰/۶۱	۰/۶۱۰۱
۲۶	Tetrahydrofuran	O1CCCC1	۰/۴۳۴۱	۰/۳۹۹۱
۲۷ <sup>t</sup>	1,4-Dioxane	O1CCOCC1	۰/۴۷۶۶	۰/۴۷۳۵
۲۸ <sup>t</sup>	Methyl phenyl ether	O(c1ccccc1)C	۰/۶۴۶۶	۰/۵۷۷۵
۲۹	Phenetole	CCOc1ccccc1	۰/۶۹۶۲	۰/۶۶۱۸
۳۰	Nitroethane	CCN(=O)=O	۰/۴۰۲۵	۰/۴۱۵۲
۳۱ <sup>t</sup>	1-Nitropropane	C(CC)N(=O)=O	۰/۴۷۱۱	۰/۴۸۴۳
۳۲ <sup>v</sup>	1-Nitrobutane	C(CCC)N(=O)=O	۰/۵۴۵۱	۰/۵۵۹۵
۳۳	1-Nitropentane	C(CCCC)N(=O)=O	۰/۶۱۷۷	۰/۵۹۷۵
۳۴	1-Nitrobenzene	c1(ccccc1)N(=O)=O	۰/۷۴۶۶	۰/۷۴۶۳
۳۵	1-Cyanoethane	C(C)C#N	۰/۳۴۵۹	۰/۳۲۴۵
۳۶	1-Cyanopropane	C(CCC)C#N	۰/۴۱۷۴	۰/۴۰۶۷
۳۷	1-Cyanobutane	C(CCCC)C#N	۰/۴۹۴	۰/۴۶۵۸
۳۸	1-Cyanopentane	C(CCCCC)C#N	۰/۵۶۶۹	۰/۵۲۷۵
۳۹	Pyridine	c1cccn1	۰/۵۱۳۸	۰/۵۱۳۹
۴۰	2-Picoline	c1(cccn1)C	۰/۵۷۰۵	۰/۵۶۵۷
۴۱ <sup>v</sup>	3-Picoline	c1c(ccn1)C	۰/۶۰۰۳	۰/۵۷۵۸
۴۲	4-Picoline	c1cc(cen1)C	۰/۵۹۸۹	۰/۵۷۴۶
۴۳	2,3-Lutidine	c1(c(cccn1)C)C	۰/۶۶۴۷	۰/۶۲۲۵
۴۴	2,4-Lutidine	c1(cc(ccn1)C)C	۰/۶۵۳۹	۰/۶۱۴۸
۴۵ <sup>v</sup>	2,5-Lutidine	c1(ccc(en1)C)C	۰/۶۵۳۹	۰/۶۳۱۴
۴۶	2,6-Lutidine	c1(cccc(n1)C)C	۰/۶۲۰۹	۰/۶۰۸۴
۴۷ <sup>t</sup>	3,4-Lutidine	c1c(c(ccn1)C)C	۰/۷۰۰۷	۰/۶۴۰۶
۴۸	3,5-Lutidine	c1c(cc(en1)C)C	۰/۶۸۵۴	۰/۶۲۸۲
۴۹	3-Chloropyridine	c1c(cccn1)Cl	۰/۶۲۹۳	۰/۶۲۳۹

## ادامه جدول ۲-۲۳

ردیف	ساختار شیمیایی	SMILES	RI واقعی	RI پیش بینی شده
۵۰	1-Acetoxyp propane	C(CC)OC(=O)C	۰/۴۵۴۲	۰/۴۷۹
۵۱	1-Acetoxyp butane (butyl acetate)	C(CCC)OC(=O)C	۰/۵۲۷۳	۰/۵۴۱
۵۲	1-Acetoxyp entane (pentyl acetate)	C(CCCC)OC(=O)C	۰/۵۹۹۴	۰/۶۱۷۱
۵۳	1,1,1-Trifluorooctane	C(CCCCCC)(F)(F)F	۰/۵۱۸۹	۰/۵۳۴۹
۵۴ <sup>v</sup>	Fluorobenzene	c1(ccccc1)F	۰/۴۷۴۲	۰/۴۸۴۴
۵۵	Hexafluorobenzene	c1(c(c(c(c(c1F)F)F)F)F)F	۰/۳۹۲۸	۰/۳۸۹۸
۵۶	Trifluoromethylbenzene	c1(ccccc1)C(F)(F)F	۰/۴۶۹۱	۰/۴۷۳۴
۵۷	Dichloromethane	C(Cl)Cl	۰/۳۵۵	۰/۳۵۴۹
۵۸	Trichloromethane	C(Cl)(Cl)Cl	۰/۴۲۹۷	۰/۴۵۹۵
۵۹ <sup>t</sup>	Tetrachloromethane	C(Cl)(Cl)(Cl)Cl	۰/۴۸۲۳	۰/۵۶۲۳
۶۰	1-Chlorobutane	C(CCC)Cl	۰/۴۵۵۹	۰/۴۴۲۶
۶۱ <sup>t</sup>	1-Chloropentane	C(CCCC)Cl	۰/۵۲۹۱	۰/۵۲۰۱
۶۲	1-Chlorohexane	C(CCCCC)Cl	۰/۶۰۱۷	۰/۵۷۱۴
۶۳ <sup>v</sup>	Chlorobenzene	c1(ccccc1)Cl	۰/۶۱۵۳	۰/۵۹۲
۶۴	1-Bromopropane	c1(ccccc1)Cl	۰/۴۴۷۵	۰/۵۸۰۸
۶۵	1-Bromobutane	C(CCC)Br	۰/۵۲۱۳	۰/۴۵۷۵
۶۶ <sup>v</sup>	1-Bromopentane	C(CCCC)Br	۰/۵۹۴۶	۰/۵۴۷۲
۶۷ <sup>v</sup>	Bromobenzene	c1(ccccc1)Br	۰/۶۸۲۸	۰/۶۶۲۹
۶۸ <sup>t</sup>	1-Butanethiol	C(CCC)S	۰/۵۱۰۶	۰/۴۷۵
۶۹ <sup>v</sup>	1-Pentanethiol	C(CCCC)S	۰/۵۸۳۷	۰/۵۶۳۳
۷۰	n-Hexanethiol	C(CCCCC)S	۰/۶۵۶۴	۰/۶۴۰۲
۷۱ <sup>t</sup>	Thiophene	c1cccs1	۰/۴۸۵۹	۰/۵۱۱
۷۲ <sup>v</sup>	1-Hexene	C=CCCCC	۰/۴۱۷۴	۰/۴۴۵
۷۳ <sup>t</sup>	Cyclohexene	C1=CCCCC1	۰/۵۰۳	۰/۵۰۹۴
۷۴	1,4-Cyclohexadiene	C1=CCC=CC1	۰/۵۱۶۱	۰/۵۰۵۴
۷۵	1,3-Cyclohexadiene	C1=CC=CCC1	۰/۴۹۱۳	۰/۵۰۳۸
۷۶	1-Heptene	C=CCCCCC	۰/۴۸۹۱	۰/۵۰۲۲
۷۶ <sup>v</sup>	1-Octene	C=CCCCCCC	۰/۵۵۹۹	۰/۵۷۳۹
۷۸ <sup>t</sup>	1-Nonene	C=CCCCCCCC	۰/۶۳۱۹	۰/۶۴۵۱
۷۹ <sup>v</sup>	1-Decene	C=CCCCCCCCC	۰/۷۰۳۱	۰/۷۰۹
۸۰	1-Pentyne	C#CCCC	۰/۳۴۴	۰/۳۹
۸۱	1-Hexyne	C#CCCCC	۰/۴۱۹۱	۰/۴۴۱۸
۸۲ <sup>t</sup>	2-Hexyne	CC#CCCC	۰/۴۵۹۶	۰/۴۶۰۹
۸۳ <sup>v</sup>	3-Hexyne	CCC#CCC	۰/۴۴۳۱	۰/۴۴۵۱
۸۴ <sup>v</sup>	1-Heptyne	C#CCCCC	۰/۴۹۱۳	۰/۵۰۶۳
۸۵	1-Octyne	C#CCCCCCC	۰/۵۶۲۴	۰/۵۷۱۷
۸۶ <sup>v</sup>	4-Octyne	CCCC#CCCC	۰/۵۷۸۷	۰/۵۹۲۳
۸۷	1-Nonyne	C#CCCCCCCC	۰/۶۳۴۱	۰/۶۴۵۵
۸۸	1-Decyne	C#CCCCCCCCC	۰/۷۰۵۵	۰/۷۱۰۵
۸۹ <sup>v</sup>	Benzene	c1ccccc1	۰/۴۸۱۱	۰/۴۷۷۳
۹۰	Toluene	c1(ccccc1)C	۰/۵۵۸۶	۰/۵۳۶۸
۹۱ <sup>v</sup>	Ethylbenzene	c1(ccccc1)CC	۰/۶۲۲۹	۰/۶۱۳
۹۲ <sup>v</sup>	Naphthalene	c1cccc2ccccc12	۰/۸۶۳	۰/۷۷۲۵

ادامه جدول ۲-۲۳

ردیف	ساختار شیمیایی	SMILES	RI واقعی	RI پیش بینی شده
۹۳ <sup>t</sup>	Azulene	c1ccc2ccccc12	۰/۹۴۲۸	۰/۷۸۹۷
۹۴	Pentane	CCCCC	۰/۳۵۷۱	۰/۳۷۴۳
۹۵	Cyclopentane	C1CCCC1	۰/۴۱۹۹	۰/۴۳
۹۶	2,2-Dimethylbutane	CC(CC)(C)C	۰/۳۸۶۵	۰/۳۹۱۳
۹۷	2,3-Dimethylbutane	CC(C(C)C)C	۰/۴۰۹	۰/۴۰۰۲
۹۸ <sup>v</sup>	Hexane	CCCCCC	۰/۴۲۸۶	۰/۴۵۱۱
۹۹ <sup>t</sup>	Cyclohexane	C1CCCCC1	۰/۴۹۵۱	۰/۵۱۱۷
۱۰۰ <sup>v</sup>	2,2-Dimethylpentane	CC(CCC)(C)C	۰/۴۴۸۹	۰/۴۵۳۸
۱۰۱	2,3-Dimethylpentane	CC([C@@H](CC)C)C	۰/۴۸۵۲	۰/۴۷۹۳
۱۰۲	2,4-Dimethylpentane	CC(CC(C)C)C	۰/۴۴۸۹	۰/۴۷۰۶
۱۰۳	2,2,3-Trimethylbutane	CC(C(C)C)(C)C	۰/۴۶۴۶	۰/۴۳۳۷
۱۰۴ <sup>t</sup>	Heptane	CCCCCCC	۰/۵	۰/۵۱۶۸
۱۰۵	Cycloheptane	C1CCCCC1	۰/۵۹۷	۰/۵۹۵۵
۱۰۶	Methylcyclohexane	C1(CCCCC1)C	۰/۵۴۱۱	۰/۵۷۰۵
۱۰۷	2,3-Dimethylhexane	CC([C@@H](CCC)C)C	۰/۵۴۷	۰/۵۳۴
۱۰۸ <sup>t</sup>	2,4-Dimethylhexane	CC(C[C@@H](CC)C)C	۰/۵۲۳۷	۰/۵۲۲۹
۱۰۹ <sup>t</sup>	3,4-Dimethylhexane	CC[C@@H]([C@@H](CC)C)C	۰/۵۵۶۶	۰/۵۳۴۸
۱۱۰	2,2,4-Trimethylpentane	CC(CC(C)C)(C)C	۰/۴۹۷۵	۰/۴۸۷۲
۱۱۱ <sup>v</sup>	2,3,4-Trimethylpentane	CC(C(C)C)C	۰/۵۴۷	۰/۴۹۸۶
۱۱۲	cis-1,2-Dimethylcyclohexane	[C@@H]1([C@H](CCCC1)C)C	۰/۶۱۸۹	۰/۶۲۲۶
۱۱۳	trans-1,2-Dimethylcyclohexane	[C@@H]1([C@@H](CCCC1)C)C	۰/۵۹۷۶	۰/۶۲۲۵
۱۱۴	cis-1,4-Dimethylcyclohexane	[C@@H]1(CC[C@H](CC1)C)C	۰/۵۹۷۴	۰/۶۲۱۴
۱۱۵ <sup>t</sup>	trans-1,4-Dimethylcyclohexane	[C@@H]1(CC[C@@H](CC1)C)C	۰/۵۸۲	۰/۶۱۸۹
۱۱۶	Octane	CCCCCCCC	۰/۵۷۱۴	۰/۵۸۰۹
۱۱۷	Cyclooctane	C1CCCCCCC1	۰/۶۸۷۹	۰/۶۶۶۷
۱۱۸ <sup>t</sup>	Nonane	CCCCCCCCC	۰/۶۴۲۹	۰/۶۵۶۳
۱۱۹	Decane	CCCCCCCCCC	۰/۷۱۴۳	۰/۷۱۱۹
۱۲۰	Cyclodecane	C1CCCCCCCC1	۰/۸۴۰۹	۰/۸۲۱۱
۱۲۱ <sup>t</sup>	cis-Hydrindane	C1CCc2ccccc12	۰/۷۳۷۱	۰/۷۲۰۱
۱۲۲	trans-Hydrindane	C1CCc2ccccc12	۰/۷۱۳۶	۰/۷۱۰۵
۱۲۳	cis-Decalin	C1CCC[C@H]2CCCC[C@@H]12	۰/۸۲۱۸	۰/۸۲۳۴
۱۲۴	trans-Decalin	C1CCC[C@H]2CCCC[C@H]12	۰/۷۹۲۲	۰/۸۱۹۶
۱۲۵ <sup>t</sup>	Adamantane	[C@@H]12C[C@@H]3C[C@@H](C1)C[C@@H]2C3	۰/۸۰۹۷	۰/۸۱۶۶
۱۲۶ <sup>t</sup>	Undecane	CCCCCCCCCCC	۰/۷۸۵۷	۰/۷۹۴۶
۱۲۷	Dodecane	CCCCCCCCCCCC	۰/۸۵۷۱	۰/۸۴۸۶
۱۲۸	Tridecane	CCCCCCCCCCCCC	۰/۹۲۸۶	۰/۹۲۷۹
۱۲۹	Tetradecane	CCCCCCCCCCCCCC	۱/۰۰۰	۰/۹۹۵۹
۱۳۰ <sup>t</sup>	Tetramethylsilane	[Si](C)(C)(C)C	۰/۳۰۳۹	۰/۳۳۰۹
۱۳۱	Hexamethyldisilane	[Si]([Si](C)(C)C)(C)C	۰/۴۹۱۵	۰/۴۹۰۴
۱۳۲	Tetramethyltin	[Sn](C)(C)(C)C	۰/۴۳۰۴	۰/۴۲۹۹

جدول ۲-۲۴ ساختار شیمیایی و SMILES مربوط به ترکیبات مورد مطالعه مجموعه داده‌های B به همراه شاخص بازداری

ردیف	ساختار شیمیایی	SMILES	RI واقعی	RI پیش بینی شده
۱ <sup>v</sup>	2-Methyl-1H-pyrrole	<chem>c1(cc(cc1C(C)(C)C)C(C)(C)C)O</chem>	۸۲۱	۸۵۶
۲ <sup>v</sup>	2-Ethyl-2,5-dimethylpyrazine	<chem>CCCCCCCCCCCC</chem>	۱۰۹۱	۱۰۸۷
۳	Vinylpyrazine	<chem>C([C@H](C(=O)OC[C@H](CC(C)C)C)C)O</chem>	۹۶۵	۹۴۷
۴	2-Furanmethanol	<chem>C(=O)/C=C/CCCCCCC</chem>	۹۱۰	۸۵۳
۵ <sup>t</sup>	Pyrazine	<chem>C(=O)/C=C/C=CCCC</chem>	۷۵۵	۷۴۰
۶	Benzaldehyde	<chem>C(=O)/C=C\C=C/CCCC</chem>	۱۰۱۴	۹۸۲
۷	Menthol	<chem>CCCCCCCCCCCC</chem>	۱۱۹۴	۱۱۹۰
۸ <sup>v</sup>	Methylvinylpyrazine	<chem>c1cccc2cccc12</chem>	۱۰۲۹	۱۰۳۸
۹	2-Pentylthiophene	<chem>CCCCCCCCCCC</chem>	۱۲۲۴	۱۱۷۰
۱۰	2-Ethyl-5-methylpyrazine	<chem>[C@H]1([C@H](CC[C@H](C1)C)C(C)C)O</chem>	۱۰۰۱	۱۰۱۵
۱۱ <sup>t</sup>	Nonanal	<chem>C(c1cccc1)n1cccc1</chem>	۱۰۸۵	۱۱۱۳
۱۲	(Z,Z)2,4-decadienal	<chem>c1(cccs1)CCCC</chem>	۱۲۲۶	۱۳۰۹
۱۳	Tridecane	<chem>c1(ccccc1)CCCC</chem>	۱۲۵۸	۱۲۹۹
۱۴ <sup>t</sup>	2-Pentylfuran	<chem>C(=O)CCCCCCCC</chem>	۱۱۲۹	۹۹۶
۱۵	2,5-Diethylpyrazine	<chem>CCCCCCCCCCC</chem>	۱۱۰۴	۱۰۹۱
۱۶ <sup>v</sup>	1-(2-furanmethyl)-1H-pyrrole	<chem>c1(enc(en1)CC)CC</chem>	۱۲۰۷	۱۱۸۲
۱۷ <sup>v</sup>	Dodecane	<chem>c1(c(nc(en1)C)CC)C</chem>	۱۱۹۲	۱۱۹۹
۱۸	3-Methylphenol	<chem>c1(c(nc(en1)C)CC)C</chem>	۱۰۱۶	۱۰۸۰
۱۹ <sup>t</sup>	3-Ethyl-2,5-dimethylpyrazine	<chem>Cc1cc(O)ccc1</chem>	۱۰۸۴	۱۰۸۵
۲۰ <sup>t</sup>	Tetradecane	<chem>c1(ccccc1)CCCC</chem>	۱۳۳۳	۱۳۹۹
۲۱	2,4-decadienal (E,E)	<chem>CC(=O)/C=C/CCCC</chem>	۱۲۴۹	۱۳۱۹
۲۲	2-Undecenal	<chem>CC1=CC[C@H](C(=C)C)CC1</chem>	۱۲۷۶	۱۳۶۵
۲۳	2-Butenal	<chem>c1(c(ncn1)C)C=C</chem>	۷۰۵	۶۴۸
۲۴	Butylated hydroxytoluene	<chem>C([C@H](CCCC)CC)O</chem>	۱۴۹۰	۱۵۱۵
۲۵ <sup>t</sup>	Propanoic acid, 2-methyl-, 3-hydroxy-2,4,4-trimethylpentyl esters	<chem>c1(ncnc1)C=C)C</chem>	۱۴۸۰	۱۳۸۴
۲۶ <sup>v</sup>	Naphthalene	<chem>c1(cnc(en1)C)CC</chem>	۱۲۰۴	۱۲۱۸
۲۷	2-Methylpyrazine	<chem>c1(cnc(en1)C)CC</chem>	۸۵۵	۸۲۸
۲۸ <sup>t</sup>	Pentyl-benzene	<chem>CCCCCCCCCCC</chem>	۱۳۳۳	۱۱۶۸
۲۹	Pentanol	<chem>o1c(ccc1)CCCC</chem>	۸۰۱	۷۷۰
۳۰ <sup>v</sup>	Furfural	<chem>C(=C\C=C(=O)CCCC)/C</chem>	۸۶۷	۸۳۸
۳۱	Methyl heptenone	<chem>S(SC)SC</chem>	۱۰۱۶	۹۹۱
۳۲ <sup>t</sup>	2,3-Pentanedione	<chem>C=C[C@H](CCCC)O</chem>	۸۰۶	۶۹۱
۳۳ <sup>t</sup>	Decane	<chem>C(=O)c1cccc1</chem>	۱۰۳۴	۱۰۰۰
۳۴	Ethylpyrazine	<chem>C(=O)/C=C/CCCC</chem>	۹۳۷	۹۲۹
۳۵	Undecane	<chem>c1(cncn1)C=C</chem>	۱۱۱۵	۱۰۹۹
۳۶ <sup>t</sup>	Butyl-benzene	<chem>c1(cncn1)CC</chem>	۱۱۳۰	۱۰۶۵
۳۷	Limonene	<chem>c1(cnc(en1)C)C</chem>	۱۰۴۰	۱۰۴۱
۳۸	3-Octen-2-one	<chem>C(=O)CCCCCC</chem>	۹۸۰	۱۰۴۳
۳۹	2-Ethyl-6-methyl pyrazine	<chem>CC(=O)CCCC</chem>	۹۹۱	۱۰۰۹
۴۰ <sup>t</sup>	2,5-Dimethylpyrazine	<chem>Cc1ccc(cc1)C</chem>	۹۰۸	۹۲۶
۴۱	Pyrrole	<chem>C(CCCCC)O</chem>	۷۴۷	۷۵۲
۴۲	2-Ethyl-1-hexanol	<chem>o1c(ccc1)CO</chem>	۱۰۱۴	۱۰۳۲
۴۳	Hexanal	<chem>[nH]1c(ccc1)C</chem>	۸۴۲	۸۰۱
۴۴	2-Heptanone	<chem>C(=O)c1cccc1</chem>	۸۸۳	۸۸۹
۴۵ <sup>v</sup>	1-Hexanol	<chem>c1(cncn1)C</chem>	۸۷۳	۸۷۲
۴۶	1-Octen-3-ol	<chem>s1cnc(c1)C</chem>	۱۰۰۸	۹۸۲
۴۷	2-Methyl-6-vinylpyrazine	<chem>C(=O)CCCC</chem>	۱۰۴۸	۱۰۳۲
۴۸	Heptanal	<chem>C(CCCC)O</chem>	۹۲۵	۹۰۷
۴۹	2-Methyl-1H-pyrrole	<chem>[nH]1cccc1</chem>	۸۵۸	۸۳۶
۵۰ <sup>v</sup>	Xylene	<chem>c1cncn1</chem>	۹۲۵	۸۸۱
۵۱ <sup>v</sup>	2(E)-Heptenal	<chem>CC(=O)C(=O)CC</chem>	۹۶۸	۹۶۸
۵۲	Dimethyl trisulfide	<chem>C(=O)/C=C/C</chem>	۹۶۷	۹۹۰

## ۲-۵-۳ رسم و بهینه‌سازی ساختار ترکیبات آلی فرار مجموعه داده‌های متفاوت

ساختارهای شیمیایی همه ترکیبات آلی فرار با استفاده از نرم‌افزار هایپرکم ترسیم شد. فرآیند بهینه‌سازی مطابق با روش کار بخش ۱-۵-۳ و با به‌کارگیری روش نیمه تجربی AM1 انجام شد و تا رسیدن به حداقل انرژی به ۰/۰۰۱ ادامه یافت. ساختارهای بهینه شده با پسوند \*.hin ذخیره شدند.

## ۲-۵-۴ استخراج توصیف‌کننده‌های ساختاری

ساختارهای بهینه ترکیبات آلی فرار هر دو مجموعه داده به‌طور مجزا در نرم‌افزار دراگون فراخوانی شدند و سپس برای هر ترکیب به تعداد ۳۲۲۴ توصیف‌کننده مولکولی در ۲۲ دسته متفاوت محاسبه شدند.

## ۲-۵-۵ پیش‌پردازش و انتخاب توصیف‌کننده‌های مؤثر

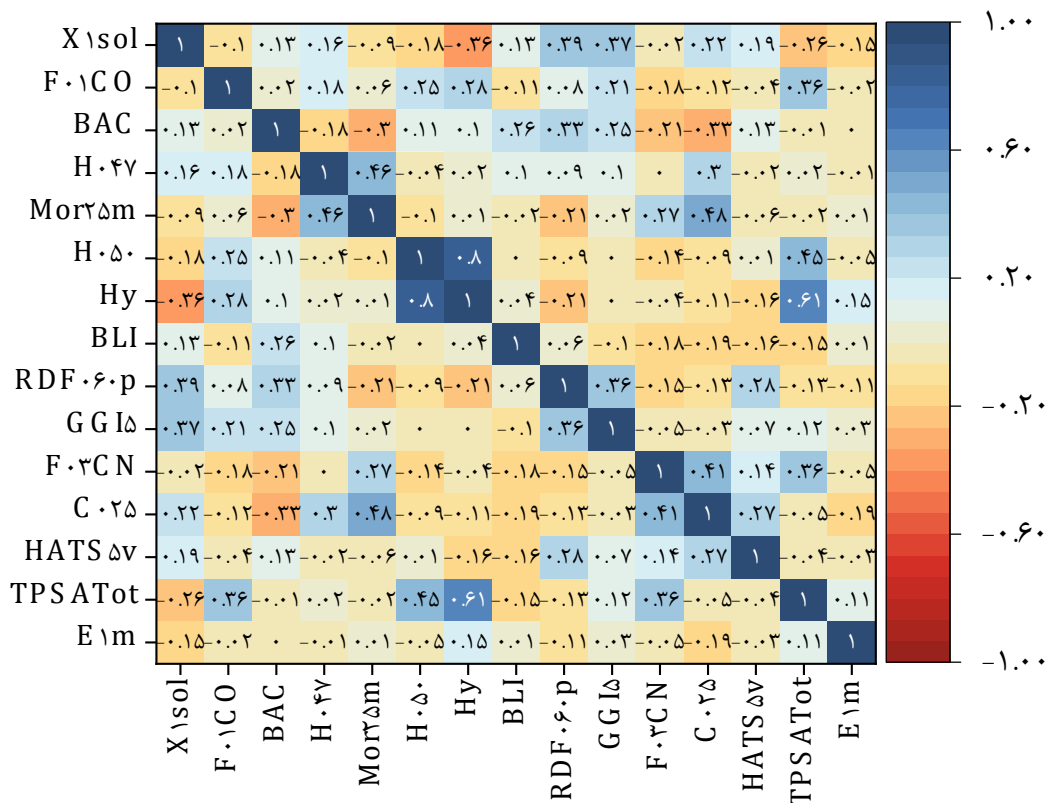
به‌منظور کاهش ابعاد داده‌ها و افزایش تفسیرپذیری مدل QSRR، باید تعداد توصیف‌کننده‌ها کاهش یابد. برای تحقق این هدف، یک استراتژی دو مرحله‌ای متشکل از مرحله غربالگری و انتخاب توصیف‌کننده انجام شد. از این‌رو، مقادیر ثابت و نسبتاً ثابت (متغیرهای با واریانس کم‌تر از ۰/۰۰۱) در مرحله غربالگری و با استفاده از اجرای بسته نرم‌افزاری caret در نرم‌افزار R حذف شدند [۱۶۳، ۱۹۵]. همچنین از بین دو توصیف‌کننده با همبستگی بالای ۰/۹، توصیف‌کننده‌ای که کم‌ترین ارتباط را با پاسخ داشت نیز حذف شد و توصیف‌کننده با بیش‌ترین همبستگی با متغیر پاسخ نگه داشته شد. پس از غربالگری، توصیف‌کننده‌های باقی‌مانده در یک ماتریس مرتب شدند، به‌طوری‌که توصیف‌کننده‌ها و RI به‌ترتیب به‌عنوان متغیرهای مستقل و وابسته در نظر گرفته شدند. داده‌های مجموعه آزمون از مجموعه داده‌ها حذف شدند. سپس، روش SCAD با استفاده از روش ارزیابی تقاطعی ده فولد (۱۰-fold-CV) موجود در بسته نرم‌افزاری ncvreg در نرم‌افزار R، روی مجموعه ارزیابی و آموزش اجرا شد [۱۶۳، ۱۹۵]. توصیف‌کننده‌های مربوط به  $\lambda_{\min}$  (با حداقل خطای ارزیابی تقاطعی (CV)) شناسایی شد و توصیف‌کننده‌های مربوط به  $\lambda_{\min}$  به‌عنوان مؤثرترین توصیف‌کننده‌ها در نظر گرفته شدند. پارامتر  $\alpha$  در محاسبات SCAD، بر اساس مطالعه

بیزی فن و لی روی ۳/۷ تنظیم شد [۳۱]. بنابراین به تعداد ۱۵ و ۸ توصیف کننده، به عنوان مؤثرترین توصیف کننده‌ها برای مجموعه داده‌های A و B به ترتیب، به دست آمد و نام و نوع توصیف کننده‌های منتخب روش SCAD برای هر دو مجموعه داده‌ها در جدول ۲-۲۵ خلاصه شدند.

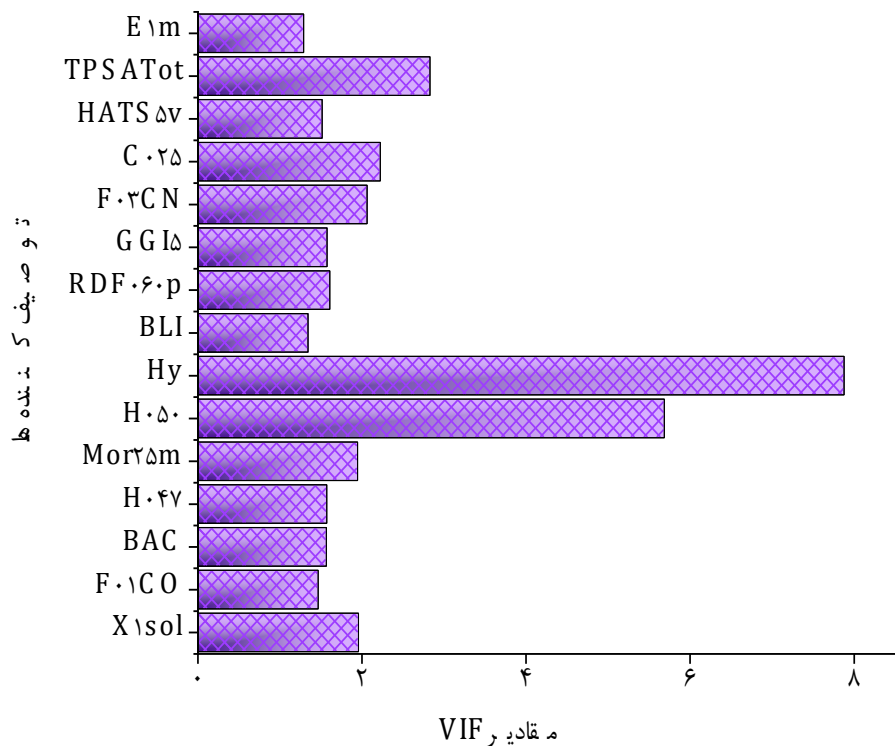
به منظور ارزیابی دقیق آماری توصیف کننده‌های منتخب روش SCAD، احتمال وجود همبستگی و هم‌خطی چندگانه نیز با محاسبه ضریب همبستگی بین دو توصیف کننده و مقادیر افزایش تورم واریانس (VIF) مورد بررسی قرار گرفتند. به این منظور نمودار نقشه رنگی و VIF برای توصیف کننده‌های منتخب روش SCAD رسم شد و نتایج حاصل در شکل ۲-۳۷ تا شکل ۲-۴۰ نمایش داده شده‌اند. نمودار نقشه رنگی هر دو مجموعه داده A و B (شکل ۲-۳۷ و شکل ۲-۳۹) نشان‌دهنده عدم وجود همبستگی معنادار بین توصیف کننده‌های منتخب روش SCAD است. با هدف بررسی پدیده هم‌خطی، پارامتر VIF محاسبه شده برای توصیف کننده‌های منتخب (شکل ۲-۳۸ و شکل ۲-۴۰) نیز نشان می‌دهد که مقادیر VIF، از مقدار هشدار ۱۰ کمتر هستند. بنابراین هم‌خطی شدید نیز بین توصیف کننده‌های منتخب روش SCAD هر دو مجموعه داده A و B وجود ندارد [۱۳۸، ۱۳۹].

جدول ۲-۲۵ توصیف‌کننده‌های منتخب SCAD برای مجموعه داده‌های A و B

مجموعه داده‌ها	ردیف	نماد	معنا	طبقه‌بندی	اهمیت
A	۱	X1sol	solvation connectivity index of order 1	Connectivity indices	۱۶/۵
	۲	F01CO	Frequency of C - O at topological distance 1	2D Frequency fingerprints	۹/۷۶
	۳	BAC	Balaban centric index	Topological indices	۸/۴۶
	۴	H-047	H attached to C1(sp3)/C0(sp2)	Atom-centred fragments	۸/۱
	۵	Mor25m	signal 25 / weighted by mass	3D-MoRSE descriptors	۷/۰۳
	۶	H-050	H attached to a heteroatom	Atom-centred fragments	۶/۷۱
	۷	Hy	hydrophilic factor	Molecular properties	۶/۶۲
	۸	BLI	Kier benzene-likeness index Topological indice	Topological indices	۶/۶۱
	۹	RDF060p	Radial Distribution Function - 060 / weighted by polarizability	RDF descriptors	۶/۰۱
	۱۰	GGI5	topological charge index of order 5	2D autocorrelations	۵/۸۷
	۱۱	F03CN	Radial Distribution Function - 030 / unweighted	2D Frequency fingerprints	۵/۳۹
	۱۲	C025	R--CR--R	Atom-centred fragments	۴/۹۴
	۱۳	HATS5v	leverage-weighted autocorrelation of lag 5 / weighted by van der Waals volume	GETAWAY descriptors	۴/۹
	۱۴	TPSA(Tot)	topological polar surface area using N,O,S,P polar contributions	Molecular properties	۴/۲۶
	۱۵	E1m	1st component accessibility directional WHIM index / weighted by mass	WHIM descriptors	۴/۲۲
B	۱	X2sol	solvation connectivity index of order 2	Connectivity indices	۱۲/۵۶
	۲	Mor07e	signal 07 / weighted by Sanderson electronegativity	3D-MoRSE descriptors	۵/۰۲
	۳	G2e	2nd component symmetry directional WHIM index / weighted by Sanderson electronegativity	WHIM descriptors	۴/۸۷
	۴	AMR	Ghose-Crippen molar refractivity	Molecular properties	۴/۵۵
	۵	TIC5	Total Information Content index (neighborhood symmetry of 5-order)	Information indices	۳/۷۱
	۶	Mor27u	signal 27 / unweighted	3D-MoRSE descriptors	۳/۱۳
	۷	TIC1	Total Information Content index (neighborhood symmetry of 1-order)	Information indices	۳/۱۱
	۸	F10CC	Frequency of C - C at topological distance 10	2D Frequency fingerprints	۱/۵۵

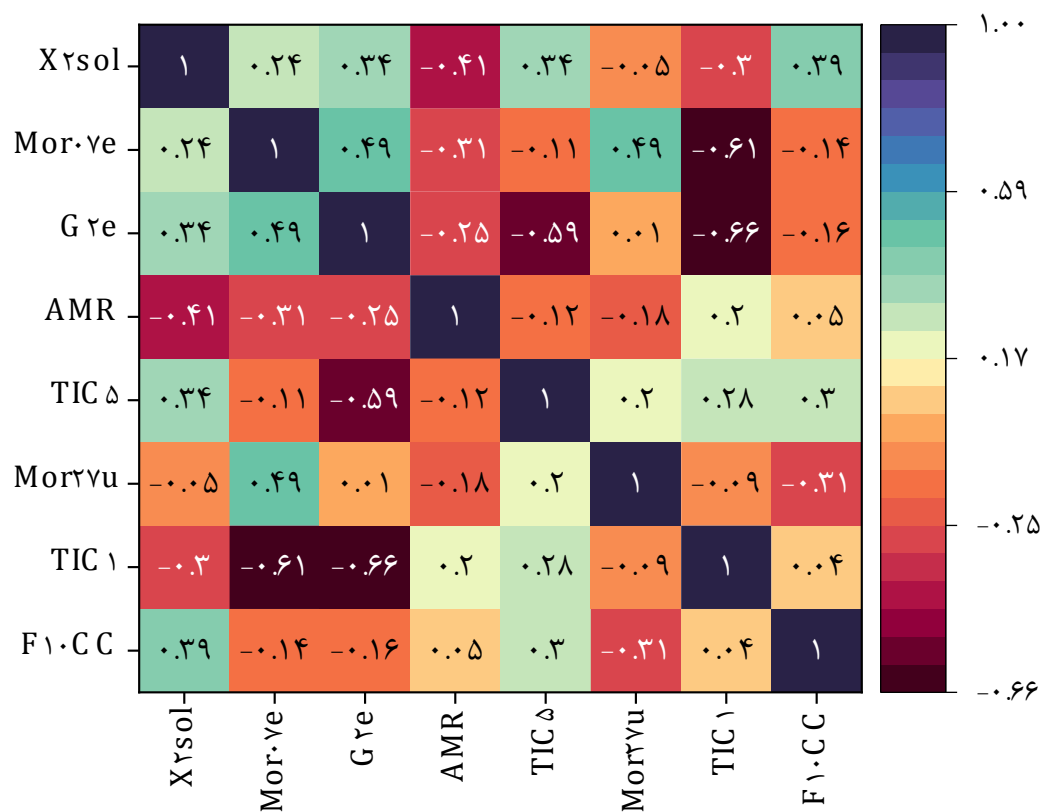


شکل ۲-۳۷ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب SCAD برای مجموعه داده A

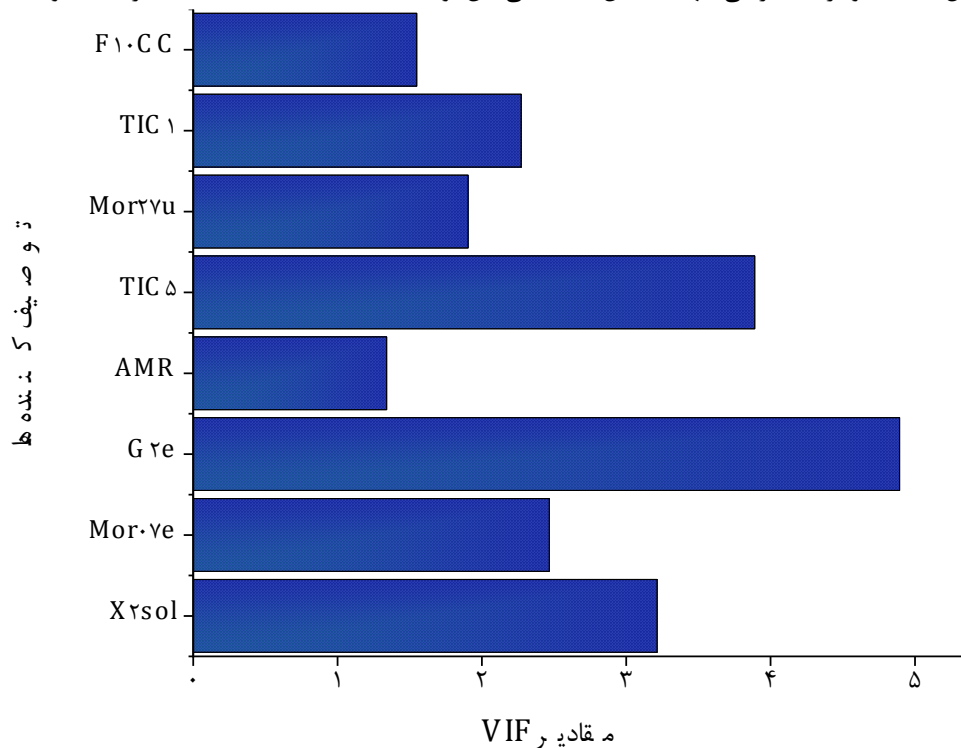


شکل ۲-۳۸ نمودار مقادیر VIF توصیف‌کننده‌های منتخب SCAD برای مجموعه داده A





شکل ۳۹-۲ نمودار نقشه رنگی جهت نمایش همبستگی بین توصیف‌کننده‌های منتخب SCAD برای مجموعه داده B



شکل ۴۰-۲ نمودار مقادیر VIF توصیف‌کننده‌های منتخب SCAD برای مجموعه داده B

## ۲-۵-۶ مدل سازی شبکه عصبی با استفاده از توصیف کننده های منتخب SCAD

مدل سازی ANN برای ایجاد یک رابطه غیرخطی بین توصیف کننده های منتخب و RI استفاده شد. در این مطالعه از یک مدل شبکه عصبی پیشخور با الگوریتم آموزشی پس انتشار خطا استفاده شد. برای آموزش مدل ANN، بهینه سازی همزمان تمام پارامترهای مهم مانند تعداد ورودی ها، تعداد گره ها در لایه پنهان، دوره های آموزشی و توابع آموزش و انتقال ضروری است [۴۰]. مدل های ANN حاوی یک لایه ورودی، یک لایه پنهان و یک لایه خروجی برای بهینه سازی پارامترهای ANN استفاده شدند. در بهینه سازی همزمان پارامترهای ANN، مدل های ANN با معماری های مختلف با استفاده از دو الگوریتم آموزشی متنوع لونبرگ مارکواریت و تنظیم بایزین (توابع آموزشی trainlm و trainbr در جعبه ابزار متلب) و دو تابع انتقال متفاوت لگاریتم سیگموئیدی و تانژانت هایپربولیک سیگموئیدی (در جعبه ابزار برنامه متلب به ترتیب با توابع logsig و tansig شناخته می شوند) برای هر دو مجموعه داده A و B به کار گرفته شد. برای بهینه سازی تعداد ورودی های شبکه عصبی مصنوعی، ابتدا تعداد نرون های لایه ورودی از ۲ تا کل توصیف کننده های منتخب SCAD در هر مجموعه داده تعریف شد. همان طور که در مطالعات قبلی گفته شد، تعداد زیر گروه هایی که می توان از تعداد توصیف کننده های منتخب تشکیل داد بسیار زیاد است و استفاده از همه این زیرگروه ها برای آموزش ANN، بسیار وقت گیر است. بنابراین از روش ANN برای یافتن ترکیب بهینه توصیف کننده های منتخب SCAD استفاده شد. مدل های ANN تولید شده با توجه به تعداد کل توصیف کننده های منتخب ساخته شدند و توابع انتقال و آموزش، گره و دور آموزشی بهینه به دست آمد. پس از آن، مقادیر همه توصیف کننده ها در محدوده تغییرات مربوط به خودشان، تصادفی شدند. هر بار مجموعه جدیدی از توصیف کننده ها با یک توصیف کننده دست کاری شده، در حضور سایر توصیف کننده های واقعی برای ساخت مدل های ANN با شرایط بهینه استفاده شد. مقادیر RI مربوط به مجموعه ارزیابی با استفاده از مدل بهینه SCAD-ANN پیش بینی شد. مقدار  $RMSE_i$  مجموعه ارزیابی با استفاده از همه توصیف کننده های منتخب و توصیف کننده

i ام دست کاری شده، محاسبه شد. این فرآیند برای همه توصیف کننده‌ها تکرار شد تا در نهایت ۱۵ و ۸ مقدار RMSE با استفاده از توصیف کننده‌های هر دو مجموعه داده A و B به دست آمد. مقدار  $RMSE_i$  بالاتر نشان می‌دهد که مدل ANN در حضور توصیف کننده i ام دست کاری شده (غیاب توصیف کننده با مقادیر واقعی خودش) دچار خطای بیش تری می‌شود. بنابراین توصیف کننده i ام با خطای بیش تر نسبت به سایر توصیف کننده‌ها از اهمیت بیش تری برخوردار است. مقادیر  $RMSE_i$  توصیف کننده‌ها به‌طور کاهشی (از نظر مقادیر RMSE) مرتب شدند و به‌عنوان ورودی ANN تعریف شدند (پارامتر اهمیت در جدول ۲-۲۵). به‌طوری که شبکه عصبی مصنوعی ابتدا با دو توصیف کننده اول با اهمیت بیش تر آموزش دید و تا تعداد کل توصیف کننده‌های مهم بهینه‌سازی ادامه یافت. به این معنی که، در بهینه‌سازی شرایط ANN، تعداد ۱۴ زیر مجموعه برای مجموعه داده A و ۷ زیر مجموعه برای مجموعه داده B با توجه به اهمیت توصیف کننده‌ها به‌عنوان ورودی استفاده شدند. تعداد نورون‌های لایه پنهان در محدوده ۲ تا ۱۰ با گام ۱ و تعداد دورهای آموزشی از ۵ تا ۵۰ (با گام ۵) تغییر یافت. مدل‌های ANN طراحی شده، با استفاده از داده‌های مجموعه آموزشی، آموزش داده شدند. در نتیجه مدل‌های توسعه یافته برای پیش‌بینی مقادیر RI مجموعه ارزیابی استفاده شد. بهترین مدل ANN با توجه به کم‌ترین مقدار RMSE مجموعه ارزیابی انتخاب شد. بهترین معماری از مدل‌های شبکه عصبی مصنوعی با توابع آموزش و انتقال متفاوت با توجه به حداقل مقادیر RMSE در جدول ۲-۲۶ خلاصه شده است. نتایج به دست آمده نشان می‌دهد که مدل‌های شبکه عصبی مصنوعی SCAD-LM-ANN با معماری ۱-۲-۱۰ برای مجموعه داده A و SCAD-BR-ANN با معماری ۱-۴-۷ برای مجموعه داده B دارای RMSE مجموعه ارزیابی حداقل بوده و به‌عنوان مدل‌های برتر برای پیش‌بینی RI ترکیبات مورد مطالعه برگزیده شدند.

جدول ۲-۲۶ ساختارهای شبکه‌های توسعه یافته با توصیف‌کننده‌های منتخب SCAD با کمترین MSE مجموعه ارزیابی هردو مجموعه داده A و B

مجموعه داده‌ها	تعداد توصیف کننده	تابع آموزش	تابع انتقال	تعداد گره	تعداد دور آموزش	RMSE <sub>validation</sub>	R <sup>2</sup> <sub>validation</sub>
مجموعه A	۱۵	تنظیم بایزین	لگاریتم-سیگموئید	۷	۱۵	۰/۰۳	۰/۹۴
	۱۰	لونبرگ-مارکوارت	لگاریتم-سیگموئید	۲	۲۰	۰/۰۲	۰/۹۵
	۱۵	تنظیم بایزین	تانژانت-سیگموئید	۵	۱۰	۰/۰۳	۰/۹۵
	۱۰	لونبرگ-مارکوارت	تانژانت-سیگموئید	۲	۲۰	۰/۰۳	۰/۹۶
مجموعه B	۷	تنظیم بایزین	لگاریتم-سیگموئید	۴	۱۰	۲۲/۳۹	۰/۹۸
	۸	لونبرگ-مارکوارت	لگاریتم-سیگموئید	۲	۱۵	۲۶/۱۵	۰/۹۷
	۶	تنظیم بایزین	تانژانت-سیگموئید	۸	۱۵	۲۴/۷۶	۰/۹۷
	۶	لونبرگ-مارکوارت	تانژانت-سیگموئید	۲	۲۰	۳۰/۴	۰/۹۷

علاوه بر این، کارایی روش SCAD به‌عنوان روش انتخاب متغیر جریمه‌ای با روش انتخاب متغیر کلاسیک SR مقایسه شد. در این راستا، متغیرهای انتخاب شده با روش‌های SR برای هر دو مجموعه داده به‌طور مجزا و با چینش شبکه عصبی به‌عنوان ورودی ANN تعریف شدند. پس از آموزش و بهینه‌سازی پارامترهای شبکه عصبی، مقادیر RI داده‌های مجموعه آزمون با استفاده از مدل ANN بهینه (SR-LM-ANN) برای مجموعه داده A و SR-BR-ANN ۱-۲-۷ برای مجموعه داده B پیش‌بینی شدند.

## ۲-۵-۷ ارزیابی مدل SCAD-ANN

قدرت پیش‌بینی مدل‌های پیشنهادی SCAD-ANN با استفاده از رویکردهای مختلف ارزیابی شد. از جمله می‌توان به پیش‌بینی مقادیر RI مجموعه آزمون با استفاده از مدل‌های بهینه، پیش‌بینی RI کل

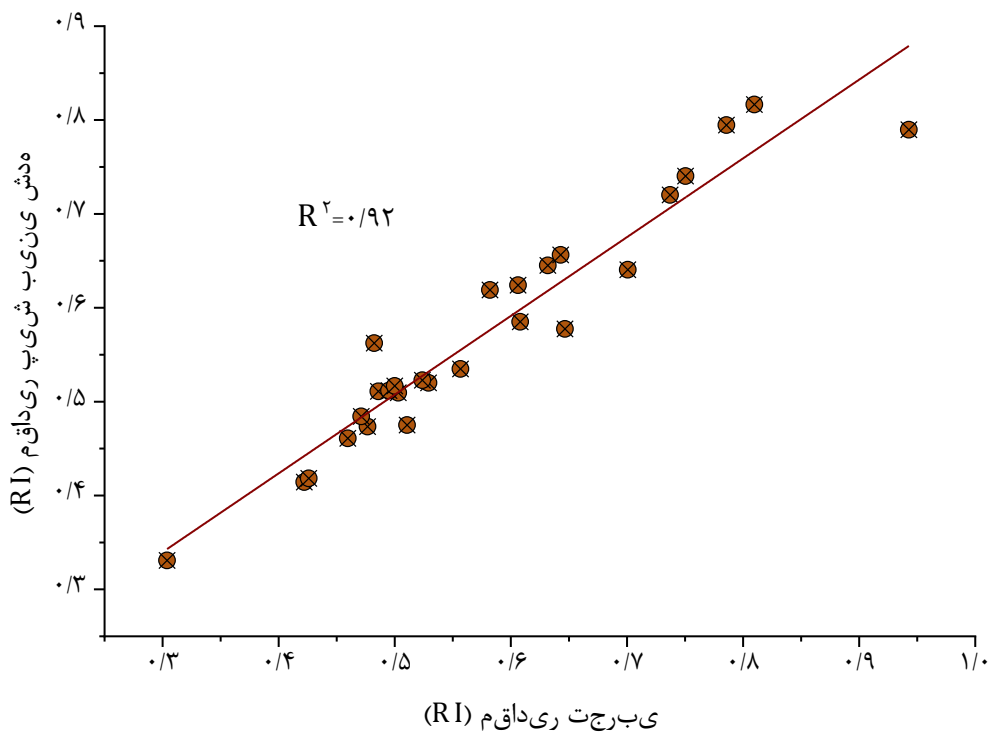
داده‌ها با استفاده از تکنیک LOO، محاسبه پارامترهای آماری، آزمون‌های پیر کاربرد دامنه کاربرد و  $Y$ -تصادفی اشاره کرد.

## ۲-۵-۷-۱ ارزیابی مدل SCAD-ANN با استفاده از پیش‌بینی مجموعه داده‌های آزمون

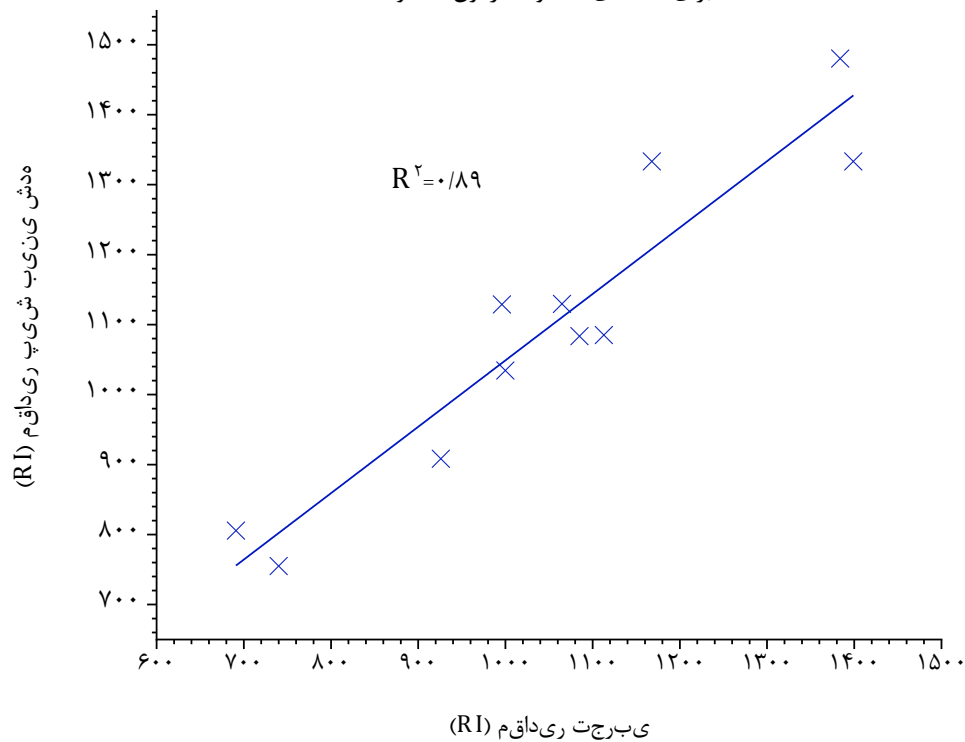
به‌منظور بررسی اعتبار و قدرت پیش‌بینی مدل، مقادیر RI مجموعه آزمون با استفاده از مدل‌های بهینه (SCAD-LM\_ANN با معماری ۱-۲-۱۰ برای مجموعه داده A و مدل SCAD-BR-ANN با معماری ۱-۴-۷ برای مجموعه داده B) پیش‌بینی شد و نتایج هر کدام در جدول ۲-۲۷ آورده شد. همان‌طور که مشاهده می‌شود مقادیر خطای پیش‌بینی اغلب ترکیبات مجموعه آزمون کم است. نتایج نشان دهنده قدرت پیش‌بینی مدل توسعه یافته است. علاوه بر این نمودار مقادیر RI پیش‌بینی شده در برابر مقادیر واقعی متناظر آن‌ها رسم شد. مقدار ضریب تعیین بالاتر از حد هشدار ( $R^2=0/6$ ) نشان‌دهنده قدرت پیش‌بینی قابل قبول مدل SCAD-ANN در هر دو مجموعه داده مورد مطالعه می‌باشد.

جدول ۲۷-۲ نتایج حاصل از ارزیابی مدل SCAD-ANN با استفاده از مجموعه آزمون

مجموعه داده‌ها	شماره ترکیبات مجموعه آزمون	شاخص بازداری (RI)		درصد خطا
		مقدار واقعی	مقدار پیش‌بینی شده	
مجموعه A	۱	۰/۴۲۲۱	۰/۴۱۴۳	-۱/۸۵
	۴	۰/۴۲۵۹	۰/۴۱۸۳	-۱/۷۹
	۶	۰/۶۰۶۲	۰/۶۲۴۱	۲/۹۶
	۱۴	۰/۷۵۰۴	۰/۷۴۰۴	-۱/۳۴
	۲۲	۰/۶۰۸	۰/۵۸۵	-۳/۷۸
	۲۷	۰/۴۷۶۶	۰/۴۷۳۵	-۰/۶۶
	۲۸	۰/۶۴۶۶	۰/۵۷۷۵	-۱۰/۶۹
	۳۱	۰/۴۷۱۱	۰/۴۸۴۳	۲/۸۱
	۴۷	۰/۷۰۰۷	۰/۶۴۰۶	-۸/۵۸
	۵۹	۰/۴۸۲۳	۰/۵۶۲۳	۱۶/۶
	۶۱	۰/۵۲۹۱	۰/۵۲۰۱	-۱/۷
	۶۸	۰/۵۱۰۶	۰/۴۷۵	-۶/۹۷
	۷۱	۰/۴۸۵۹	۰/۵۱۱	۵/۱۶
	۷۳	۰/۵۰۳	۰/۵۰۹۴	۱/۲۷
	۷۸	۰/۶۳۱۹	۰/۶۴۵۱	۲/۰۸
	۸۲	۰/۴۵۹۶	۰/۴۶۰۹	۰/۲۸
	۹۳	۰/۹۴۲۸	۰/۷۸۹۷	-۱۶/۲۴
	۹۹	۰/۴۹۵۱	۰/۵۱۱۷	۳/۳۶
	۱۰۴	۰/۵	۰/۵۱۶۸	۳/۳۶
	۱۰۸	۰/۵۲۳۷	۰/۵۲۲۹	-۰/۱۶
۱۰۹	۰/۵۵۶۶	۰/۵۳۴۸	-۳/۹۲	
۱۱۵	۰/۵۸۲	۰/۶۱۸۹	۶/۳۵	
۱۱۸	۰/۶۴۲۹	۰/۶۵۶۳	۲/۰۹	
۱۲۱	۰/۷۳۷۱	۰/۷۲۰۱	-۲/۳	
۱۲۵	۰/۸۰۹۷	۰/۸۱۶۶	۰/۸۵	
۱۲۶	۰/۷۸۵۷	۰/۷۹۴۶	۱/۱۳	
۱۳۰	۰/۳۰۳۹	۰/۳۳۰۹	۸/۸۸	
مجموعه B	۵	۷۴۰	۷۵۴/۸	۲/۰۰
	۱۱	۱۱۱۳	۱۰۸۵	-۲/۵۲
	۱۴	۹۹۶	۱۱۲۸/۸	۱۳/۳۳
	۱۹	۱۰۸۵	۱۰۸۳/۵	-۰/۱۴
	۲۰	۱۳۹۹	۱۳۳۳/۱	-۴/۷۱
	۲۵	۱۳۸۴	۱۴۸۰/۲	۶/۹۵
	۲۸	۱۱۶۸	۱۳۳۳/۱	۱۴/۱۳
	۳۲	۶۹۱	۸۰۵/۶	۱۶/۵۸
	۳۳	۱۰۰۰	۱۰۳۴/۳	۳/۴۳
۳۶	۱۰۶۵	۱۱۲۹/۶	۶/۰۶	



شکل ۴۱-۲ نمودار تغییرات مقادیر پیش‌بینی شده RI به‌وسیله مدل SCAD-LM-ANN در شرایط بهینه در مقابل مقادیر تجربی برای داده‌های مجموعه آزمون مجموعه A



شکل ۴۲-۲ نمودار تغییرات مقادیر پیش‌بینی شده RI به‌وسیله مدل SCAD-BR-ANN در شرایط بهینه در مقابل مقادیر تجربی برای داده‌های مجموعه آزمون مجموعه B

۲-۵-۷-۲ ارزیابی مدل SCAD-ANN با پیش بینی شاخص‌های بازداری تمام ترکیبات مجموعه

### داده‌های A و B با استفاده از روش رد مرحله‌ای تک تک

در این بخش، یکی از تکنیک‌های کارآمد برای ارزیابی مدل بهینه و بررسی قدرت پیش بینی مدل از تکنیک رد مرحله‌ای تک تک استفاده شد. در این روش هر ترکیب یکبار به‌عنوان داده آزمون حذف شد و مدل بهینه SCAD-ANN با استفاده از داده‌های باقی‌مانده آموزش داده شد. پس از اجرای مدل‌های بهینه برای هر دو مجموعه A و B، مقادیر RI همه ترکیبات پیش‌بینی و در جدول ۲-۲۸ و جدول ۲-۲۹ آورده شد. نتایج خطای پیش بینی اغلب ترکیبات کم بوده و این امر نشان دهنده قدرت پیش بینی مناسب مدل SCAD-ANN است. علاوه بر این، مقادیر پیش‌بینی شده بر حسب مقادیر واقعی RI برای هر مجموعه داده ترسیم شدند (شکل ۲-۴۵ و شکل ۲-۴۶ به ترتیب برای مجموعه A و B). مقادیر  $Q_{LOO}^2$  مربوط به مجموعه داده‌های A و B به ترتیب برابر با ۰/۹۴ و ۰/۸۸ به دست آمد. مقادیر  $Q_{LOO}^2$  بزرگ‌تر از ۰/۵ حاکی از اعتبار و پایداری بالای مدل توسعه یافته SCAD-ANN است.

علاوه بر این، نتایج حاصل برای هر دو مجموعه داده با استفاده از نمودارهای باقی‌مانده‌ها نیز ارزیابی شد. به این منظور، نمودارهای باقی‌مانده‌های استاندارد شده ( $\hat{\epsilon}_i$ ) بر حسب مقادیر واقعی رسم شد (شکل ۲-۴۴ و شکل ۲-۴۶). الگوی تصادفی و پراکندگی داده‌ها حول محور صفر در نمودار باقی‌مانده‌ها نشان می‌دهد که هیچ خطای سیستماتیکی در مدل‌های ANN توسعه یافته با توصیف کننده‌های منتخب روش SCAD وجود ندارد.



جدول ۲۸-۲ نتایج حاصل از ارزیابی مدل SCAD-ANN با تکنیک LOO برای کل داده‌های مجموعه A

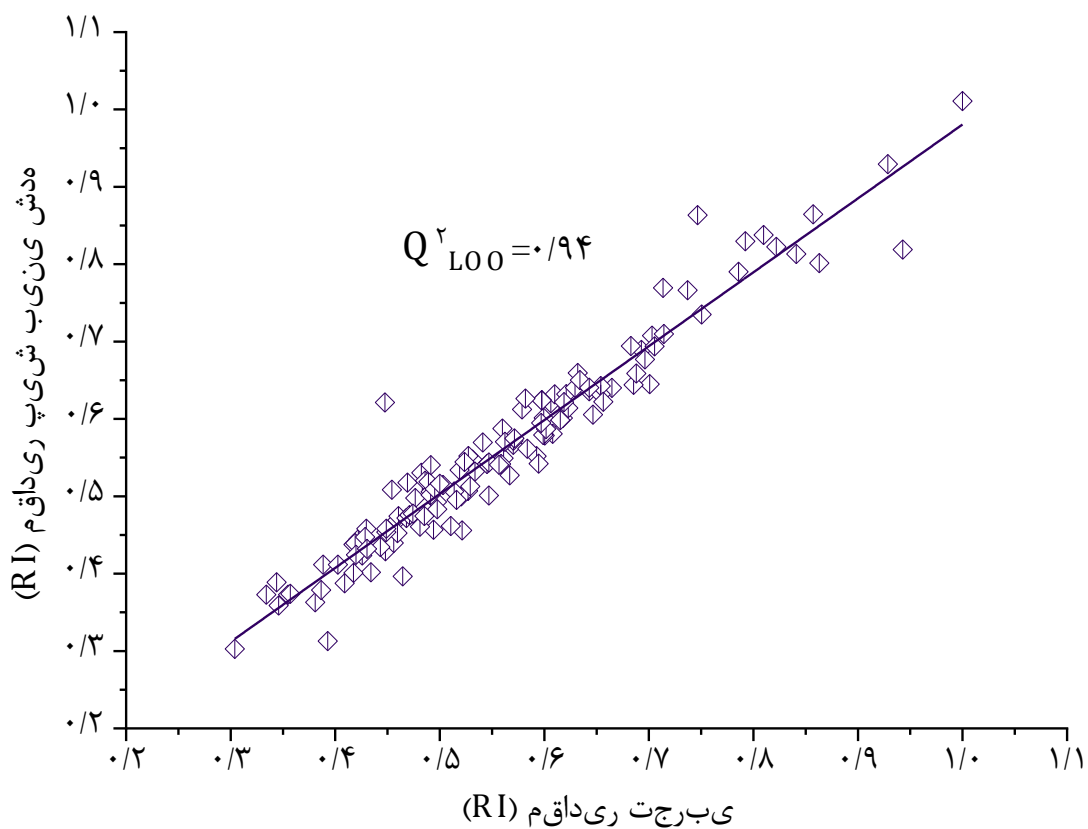
شماره ترکیب	RI			شماره ترکیب	RI		
	مقدار واقعی	مقدار پیش‌بینی شده	درصد خطا		مقدار واقعی	مقدار پیش‌بینی شده	درصد خطا
۱	۰/۴۲	۰/۴۴	-۰/۶۶	۴۰	۰/۵۷	۰/۵۷	-۰/۰۴
۲	۰/۳۳	۰/۳۷	-۱/۱۵	۴۱	۰/۱۶	۰/۵۸	۰/۶۵
۳	۰/۵	۰/۴۹	۰/۱۳	۴۲	۰/۱۶	۰/۵۸	۰/۵۸
۴	۰/۴۳	۰/۴۲	۰/۰۷	۴۳	۰/۶۶	۰/۶۴	۰/۷۳
۵	۰/۵۷	۰/۵۷	۰/۱۱	۴۴	۰/۶۵	۰/۶۳	۰/۶۶
۶	۰/۶۱	۰/۶۲	-۰/۴۳	۴۵	۰/۶۵	۰/۶۴	۰/۳۴
۷	۰/۴۹	۰/۵۱	-۰/۵۱	۴۶	۰/۶۲	۰/۶۳	-۰/۳۳
۸	۰/۶۴	۰/۶۴	۰/۲۱	۴۷	۰/۷	۰/۶۵	۱/۶۵
۹	۰/۵۶	۰/۵۵	۰/۳۵	۴۸	۰/۶۹	۰/۶۴	۱/۲۳
۱۰	۰/۳۹	۰/۴۱	-۰/۶۹	۴۹	۰/۶۳	۰/۶۴	-۰/۳۳
۱۱	۰/۴۶	۰/۴۷	-۰/۲۹	۵۰	۰/۴۵	۰/۵۱	-۱/۶۱
۱۲	۰/۵۳	۰/۵۵	-۰/۴۸	۵۱	۰/۵۳	۰/۵۵	-۰/۷۴
۱۳	۰/۱۶	۰/۵۹	۰/۳۳	۵۲	۰/۱۶	۰/۱۶	-۰/۰۶
۱۴	۰/۷۵	۰/۷۴	۰/۴۵	۵۳	۰/۵۲	۰/۵۳	-۰/۴۵
۱۵	۰/۶۹	۰/۶۹	۰/۱۱	۵۴	۰/۴۷	۰/۴۸	-۰/۱
۱۶	۰/۴۶	۰/۴۷	-۰/۴۱	۵۵	۰/۳۹	۰/۳۱	۲/۳۸
۱۷	۰/۵۳	۰/۵۳	۰/۰۵	۵۶	۰/۴۷	۰/۵۲	-۱/۴۴
۱۸	۰/۳۸	۰/۳۶	۰/۵۳	۵۷	۰/۳۶	۰/۳۷	-۰/۵۳
۱۹	۰/۴۵	۰/۴۳	۰/۵۷	۵۸	۰/۴۳	۰/۴۶	-۰/۸۵
۲۰	۰/۵۳	۰/۵۱	۰/۵۹	۵۹	۰/۴۸	۰/۵۳	-۱/۴۳
۲۱	۰/۵۲	۰/۵	۰/۶۱	۶۰	۰/۴۶	۰/۴۴	۰/۴۸
۲۲	۰/۶۱	۰/۵۸	۰/۸۱	۶۱	۰/۵۳	۰/۵۱	۰/۴۸
۲۳	۰/۵۹	۰/۵۵	۱/۲۱	۶۲	۰/۱۶	۰/۵۹	۰/۴۴
۲۴	۰/۴۷	۰/۴۷	-۰/۱۱	۶۳	۰/۶۲	۰/۱۶	۰/۵
۲۵	۰/۶۱	۰/۶۳	-۰/۶۵	۶۴	۰/۴۵	۰/۶۲	-۳/۱۶
۲۶	۰/۴۳	۰/۴	۰/۹۶	۶۵	۰/۵۲	۰/۴۶	۱/۹۴
۲۷	۰/۴۸	۰/۵	-۰/۶۳	۶۶	۰/۵۹	۰/۵۴	۱/۵۵
۲۸	۰/۶۵	۰/۶۱	۱/۲۲	۶۷	۰/۶۸	۰/۶۹	-۰/۳۴
۲۹	۰/۷	۰/۶۸	۰/۵۷	۶۸	۰/۵۱	۰/۴۶	۱/۴۶
۳۰	۰/۴	۰/۴۱	-۰/۲۷	۶۹	۰/۵۸	۰/۵۶	۰/۶۶
۳۱	۰/۴۷	۰/۴۸	-۰/۱۳	۷۰	۰/۶۶	۰/۶۲	۱/۰۲
۳۲	۰/۵۵	۰/۵۴	۰/۱	۷۱	۰/۴۹	۰/۵۲	-۱/۰۱
۳۳	۰/۶۲	۰/۱۶	۰/۴۹	۷۲	۰/۴۲	۰/۴۴	-۰/۶۱
۳۴	۰/۷۵	۰/۸۶	-۳/۴۷	۷۳	۰/۵	۰/۵۱	-۰/۳۴
۳۵	۰/۳۵	۰/۳۶	-۰/۳۸	۷۴	۰/۵۲	۰/۵	۰/۶۲
۳۶	۰/۴۲	۰/۴	۰/۴۹	۷۵	۰/۴۹	۰/۵۱	-۰/۴۳
۳۷	۰/۴۹	۰/۴۶	۱/۱۱	۷۶	۰/۴۹	۰/۵۲	-۰/۸۸
۳۸	۰/۵۷	۰/۵۳	۱/۱۹	۷۷	۰/۵۶	۰/۵۹	-۰/۸۲
۳۹	۰/۵۱	۰/۵۱	۰/۲۳	۷۸	۰/۶۳	۰/۶۶	-۰/۸۲

ادامه جدول ۲۸-۲

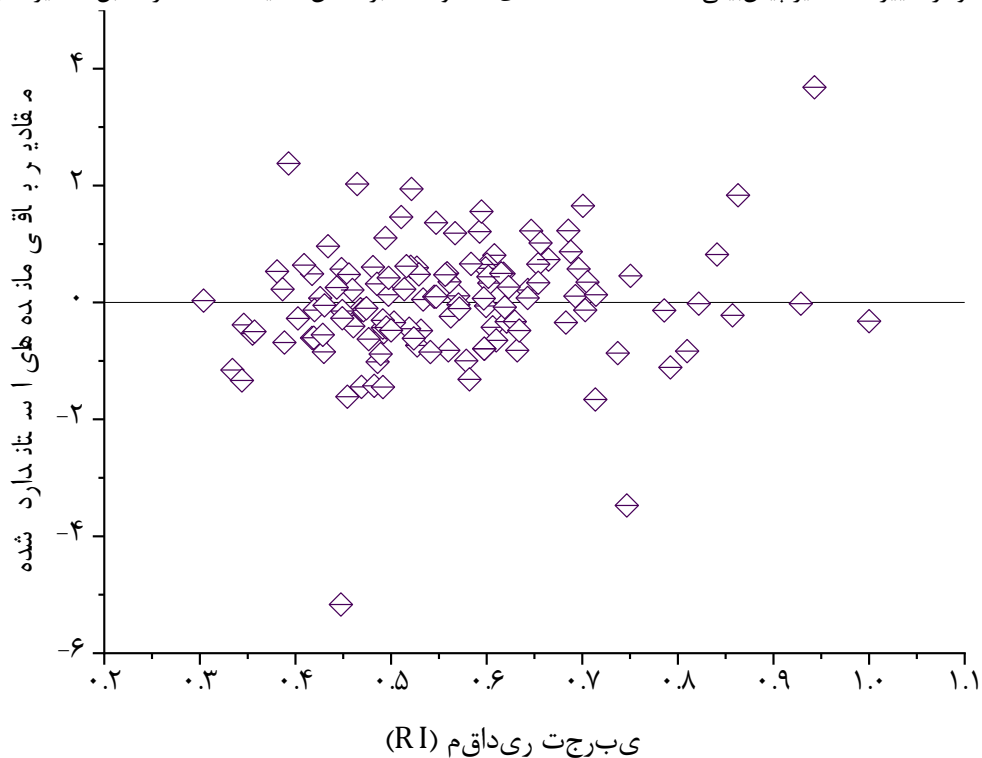
شماره ترکیب	RI		درصد خطا	شماره ترکیب	RI		درصد خطا
	مقدار پیش‌بینی شده	مقدار پیش‌بینی شده			مقدار پیش‌بینی شده	مقدار پیش‌بینی شده	
۸۱	۰/۷	۰/۷۱	-۰/۱۳	۱۰۷	۰/۵۵	۰/۵۴	۰/۱
۸۲	۰/۳۴	۰/۳۹	-۱/۳۳	۱۰۸	۰/۵۲	۰/۵۴	-۰/۶۲
۸۳	۰/۴۲	۰/۴۴	-۰/۶	۱۰۹	۰/۵۶	۰/۵۴	۰/۴۷
۸۴	۰/۴۶	۰/۴۵	۰/۲۱	۱۱۰	۰/۵	۰/۴۸	۰/۴۲
۸۵	۰/۴۴	۰/۴۳	۰/۲۶	۱۱۱	۰/۵۵	۰/۵	۱/۳۶
۸۶	۰/۴۹	۰/۵	-۰/۳۱	۱۱۲	۰/۶۲	۰/۶۲	-۰/۰۸
۸۷	۰/۵۶	۰/۵۷	-۰/۲۴	۱۱۳	۰/۶	۰/۶۲	-۰/۷۹
۸۸	۰/۵۸	۰/۶۱	۱	۱۱۴	۰/۶	۰/۶۲	-۰/۷۹
۸۹	۰/۶۳	۰/۶۵	-۰/۴۸	۱۱۵	۰/۵۸	۰/۶۳	-۱/۳۲
۹۰	۰/۷۱	۰/۶۹	۰/۳۵	۱۱۶	۰/۵۷	۰/۵۷	-۰/۱۱
۹۱	۰/۴۸	۰/۴۶	۰/۶۱	۱۱۷	۰/۶۹	۰/۶۶	۰/۸۷
۹۲	۰/۵۶	۰/۵۴	۰/۵	۱۱۸	۰/۶۴	۰/۶۴	۰/۰۸
۹۳	۰/۶۲	۰/۶۱	۰/۲۶	۱۱۹	۰/۷۱	۰/۷۱	۰/۱۳
۹۴	۰/۸۶	۰/۸	۱/۸۳	۱۲۰	۰/۸۴	۰/۸۱	۰/۸۲
۹۵	۰/۹۴	۰/۸۲	۳/۶۸	۱۲۱	۰/۷۴	۰/۷۷	-۰/۸۷
۹۶	۰/۳۶	۰/۳۷	-۰/۵	۱۲۲	۰/۷۱	۰/۷۷	-۱/۶۶
۹۷	۰/۴۲	۰/۴۲	-۰/۱۳	۱۲۳	۰/۸۲	۰/۸۲	-۰/۰۲
۹۸	۰/۳۹	۰/۳۸	۰/۲۳	۱۲۴	۰/۷۹	۰/۸۳	-۱/۱۱
۹۹	۰/۴۱	۰/۳۹	۰/۶۴	۱۲۵	۰/۸۱	۰/۸۴	-۰/۸۳
۱۰۰	۰/۴۳	۰/۴۵	-۰/۵۶	۱۲۶	۰/۷۹	۰/۷۹	-۰/۱۴
۱۰۱	۰/۵	۰/۵۱	-۰/۴۳	۱۲۷	۰/۸۶	۰/۸۶	-۰/۲۲
۱۰۲	۰/۴۵	۰/۴۵	-۰/۱۵	۱۲۸	۰/۹۳	۰/۹۳	-۰/۰۲
۱۰۳	۰/۴۹	۰/۴۷	۰/۳۲	۱۲۹	۱	۱/۰۱	-۰/۳۲
۱۰۴	۰/۴۵	۰/۴۶	-۰/۲۷	۱۳۰	۰/۳	۰/۳	۰/۰۳
۱۰۵	۰/۴۶	۰/۴	۲/۰۳	۱۳۱	۰/۴۹	۰/۵۴	-۱/۴۵
۱۰۶	۰/۵	۰/۵۲	-۰/۴۸	۱۳۲	۰/۴۳	۰/۴۳	-۰/۰۵

جدول ۲۹-۲ نتایج حاصل از ارزیابی مدل SCAD-ANN با تکنیک LOO برای کل داده‌های مجموعه B

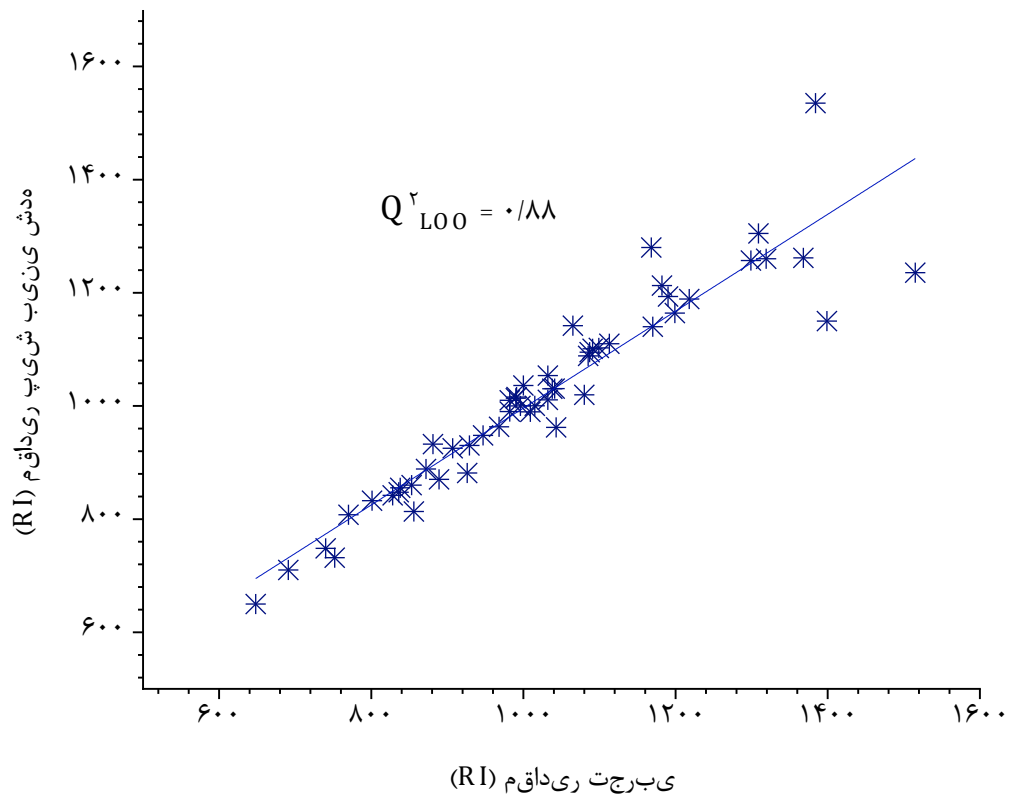
شماره ترکیب	RI			شماره ترکیب	RI		
	مقدار واقعی	مقدار پیش‌بینی شده	درصد خطا		مقدار واقعی	مقدار پیش‌بینی شده	درصد خطا
۱	۸۵۶	۸۱۳/۳۵	۴/۹۸	۲۷	۸۲۸	۸۴۱/۹۱	-۱/۶۸
۲	۱۰۸۷	۱۰۹۵/۱۴	-۰/۷۵	۲۸	۱۱۶۸	۱۲۸۰	-۹/۵۹
۳	۹۴۷	۹۴۸	-۰/۱۱	۲۹	۷۷۰	۸۰۷/۷۴	-۴/۹
۴	۸۵۳	۸۶۰	-۰/۸۲	۳۰	۸۳۸	۸۵۵/۳۶	-۲/۰۷
۵	۷۴۰	۷۴۸/۴۲	-۱/۱۴	۳۱	۹۹۱	۱۰۱۵/۶۳	-۲/۴۹
۶	۹۸۲	۹۹۰	-۰/۸۱	۳۲	۶۹۱	۷۱۰	-۲/۷۵
۷	۱۱۹۰	۱۱۹۳/۳۸	-۰/۲۸	۳۳	۱۰۰۰	۱۰۳۶/۲	-۳/۶۲
۸	۱۰۳۸	۱۰۳۰/۴۹	۰/۷۲	۳۴	۹۲۹	۹۳۰/۲۹	-۰/۱۴
۹	۱۱۷۰	۱۱۴۰	۲/۵۶	۳۵	۱۰۹۹	۱۱۰۲/۶۶	-۰/۳۳
۱۰	۱۰۱۵	۱۰۰۰/۳۷	۱/۴۴	۳۶	۱۰۶۵	۱۱۴۱/۷۷	-۷/۲۱
۱۱	۱۱۱۳	۱۱۱۰	۰/۲۷	۳۷	۱۰۴۱	۱۰۳۰/۶۹	۰/۹۹
۱۲	۱۳۰۹	۱۳۰۵	۰/۳۱	۳۸	۱۰۴۳	۹۶۲/۰۵	۷/۷۶
۱۳	۱۲۹۹	۱۲۵۶/۹۸	۳/۲۳	۳۹	۱۰۰۹	۹۸۹/۴۶	۱/۹۴
۱۴	۹۹۶	۱۰۰۰	-۰/۴	۴۰	۹۲۶	۸۸۱/۳۹	۴/۸۲
۱۵	۱۰۹۱	۱۱۰۱/۱۳	-۰/۹۳	۴۱	۷۵۲	۷۳۱/۹۱	۲/۶۷
۱۶	۱۱۸۲	۱۲۱۲/۸	-۲/۶۱	۴۲	۱۰۳۲	۱۰۱۰/۶۵	۲/۰۷
۱۷	۱۱۹۹	۱۱۶۴/۰۸	۲/۹۱	۴۳	۸۰۱	۸۳۲/۶۶	-۳/۹۵
۱۸	۱۰۸۰	۱۰۱۹/۶۷	۵/۵۹	۴۴	۸۸۹	۸۷۰/۱۷	۲/۱۲
۱۹	۱۰۸۵	۱۰۸۸/۳۹	-۰/۳۱	۴۵	۸۷۲	۸۸۸/۵۷	-۱/۹
۲۰	۱۳۹۹	۱۱۵۰	۱۷/۸	۴۶	۹۸۲	۱۰۱۰/۱	-۲/۸۶
۲۱	۱۳۱۹	۱۲۶۰/۱	۴/۴۷	۴۷	۱۰۳۲	۱۰۵۳/۵۶	-۲/۰۹
۲۲	۱۳۶۸	۱۲۶۱/۸۱	۷/۷۶	۴۸	۹۰۷	۹۲۵/۲	-۲/۰۱
۲۳	۶۴۸	۶۵۰	-۰/۳۱	۴۹	۸۳۶	۸۴۷/۴۳	-۱/۳۷
۲۴	۱۵۱۵	۱۲۳۵/۴۱	۱۸/۴۵	۵۰	۸۸۱	۹۳۲/۵۸	-۵/۸۵
۲۵	۱۳۸۴	۱۵۳۵/۲۲	-۱۰/۹۳	۵۱	۹۶۸	۹۶۳/۲۳	۰/۴۹
۲۶	۱۲۱۸	۱۱۸۹/۰۱	۲/۳۸	۵۲	۹۹۰	۱۰۱۳/۵۴	-۲/۳۸



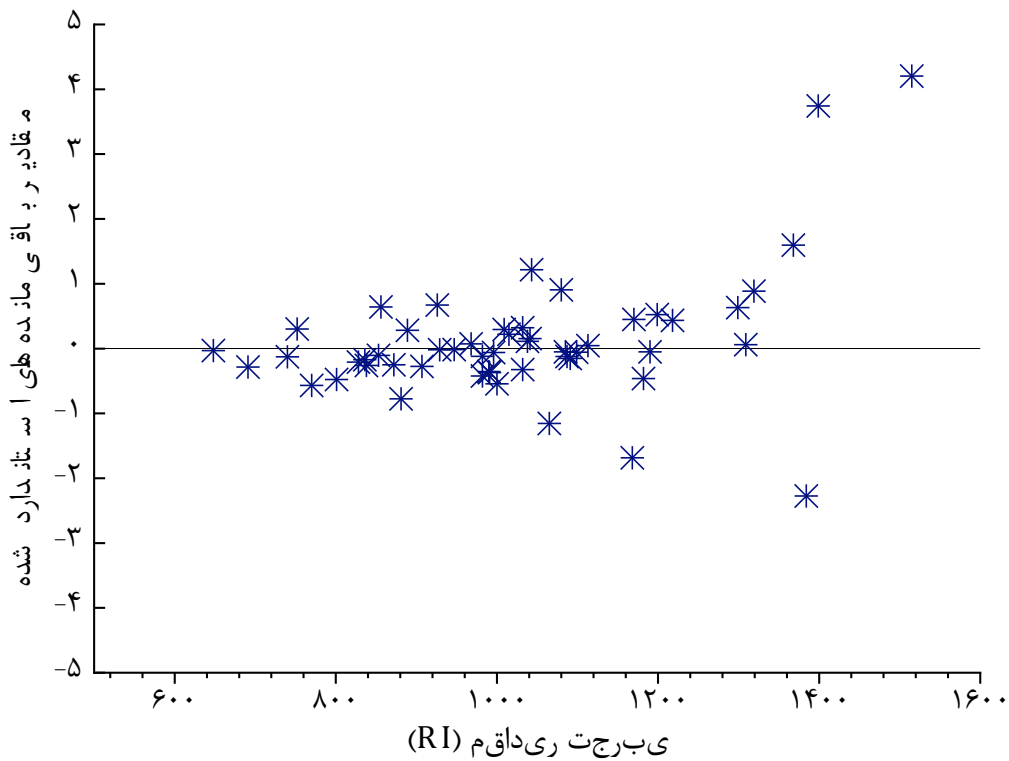
شکل ۲-۴۴ نمودار تغییرات مقادیر پیش‌بینی شده RI همه داده‌های مجموعه A بر اساس تکنیک LOO در مقابل مقادیر تجربی



شکل ۲-۴۴ نمودار باقی‌مانده‌های پیش‌بینی شده RI همه داده‌های مجموعه A با استفاده از تکنیک LOO بر حسب مقادیر تجربی



شکل ۴۵-۲ نمودار تغییرات مقادیر پیش‌بینی شده RI همه داده‌های مجموعه B بر اساس تکنیک LOO در مقابل مقادیر تجربی



شکل ۴۶-۲ نمودار باقی مانده‌های پیش‌بینی شده RI همه داده‌های مجموعه B با استفاده از تکنیک LOO بر حسب مقادیر تجربی

پارامترهای آماری  $RMSE$  و  $R^2$  مربوط به مجموعه آزمون و کل داده‌ها (تکنیک LOO) برای هر دو مجموعه داده A و B محاسبه شدند.  $RMSE$  مربوط به مجموعه آزمون برای مجموعه داده A و B به ترتیب برابر با ۰/۰۵ و ۱۰۴/۰۲ و علاوه بر این برای کل داده‌ها به ترتیب برابر با ۰/۰۵ و ۸۸/۵۰ به دست آمد. علاوه بر این ضرایب تعیین مربوط به مجموعه آزمون نیز به ترتیب برابر با ۰/۹۰ و ۰/۸۷ و ضرایب تعیین مربوط به کل داده‌ها (تکنیک LOO) به ترتیب برابر با ۰/۸۴ و ۰/۸۴ برای مجموعه داده A و B به دست آمد. با مقایسه نتایج مدل‌های برتر SCAD-ANN (جدول ۲-۲۶) برای هر دو مجموعه داده مشخص می‌شود که مدل‌های پیشنهادی در پیش‌بینی RI ترکیبات مورد نظر با قدرت بیش‌تری نسبت به مدل‌های SR-ANN عمل کرده‌اند.

## ۲-۵-۷-۳ ارزیابی مدل SCAD-ANN با استفاده از پارامترهای آماری

پارامترهای آماری مختلف نیز برای ارزیابی بیش‌تر توانایی پیش‌بینی مدل‌های SCAD-ANN پیشنهادی محاسبه شد. بنابراین پارامترهای آماری معرفی شده در بخش ۱-۵-۸-۴ برای RI پیش‌بینی شده ترکیبات مجموعه آزمون و RI پیش‌بینی شده برای کل ترکیبات به روش رد مرحله‌ای تک تک محاسبه و در جدول ۲-۳۰ خلاصه شدند. نتایج حاصل از محاسبات آماری جدول ۲-۳۰ نشان می‌دهد که پارامترهای آماری در محدوده قابل قبول قرار دارند. پارامترهای آماری تروپشا و روی (پارامترهای  $R_0^2$ ،  $R_0^2$  نسبی،  $R_m^2$  و غیره) از جمله پارامترهای آماری پر کاربرد وابسته به  $R^2$  هستند. بزرگ‌تر بودن این پارامترها از مقدار هشدار ۰/۶ و نزدیک بودن آن‌ها به  $R^2$  نشان‌دهنده تعمیم‌پذیری و قدرت پیش‌بینی رضایت بخش مدل‌های توسعه یافته QSAR مبتنی بر SCAD-ANN است. علاوه بر این شیب نمودار حاصل از مقادیر پیش‌بینی شده بر حسب مقادیر تجربی RI (و بالعکس) در عرض از مبدأ صفر نیز در محدوده ۰/۸۵ تا ۱/۱۵ قرار دارند که این نتیجه نیز حاکی از صحت مدل توسعه یافته ANN با توصیف‌کننده‌های منتخب روش SCAD می‌باشد.

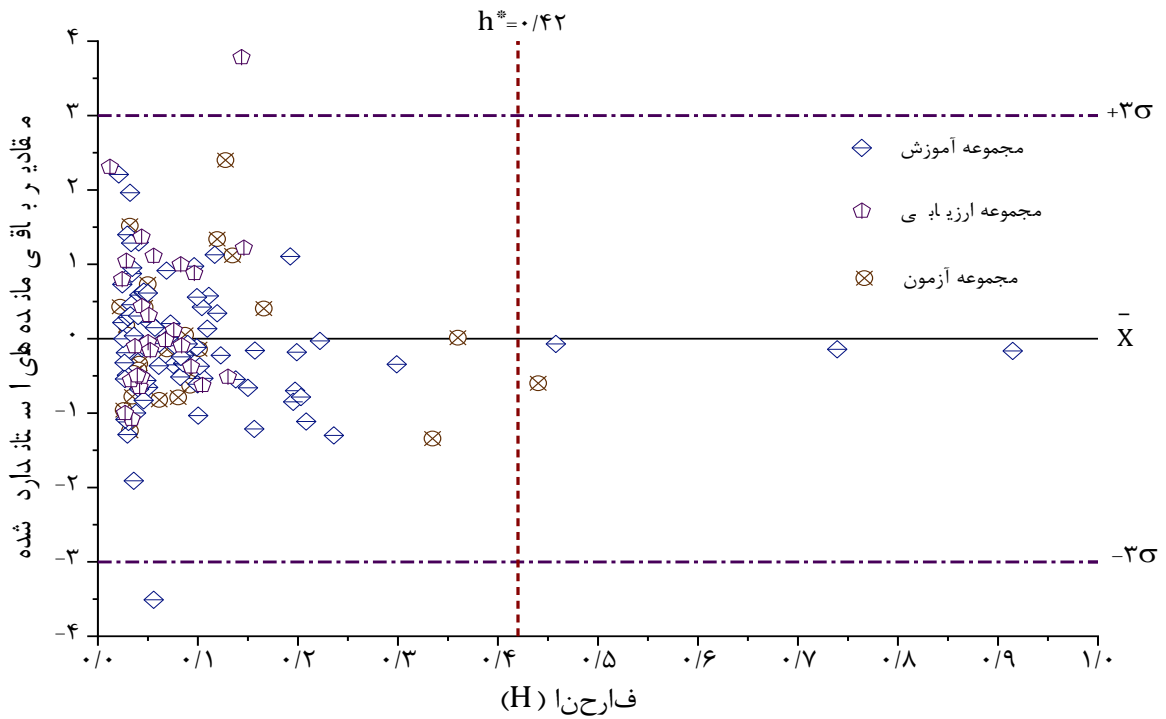
جدول ۲-۳ پارامترهای آماری محاسبه شده برای مجموعه آزمون و داده‌های پیش‌بینی شده با تکنیک LOO برای مدل SCAD-ANN هر دو مجموعه از داده‌ها

ردیف	پارامتر آماری	مجموعه داده A		مجموعه داده B		محدوده قابل قبول
		مقادیر RI پیش‌بینی شده با SCAD-ANN				
		ترکیبات مجموعه آزمون	کل ترکیبات به روش LOO	ترکیبات مجموعه آزمون	کل ترکیبات به روش LOO	
۱	PRESS	۰/۰۵	۰/۱۵	۷۸۴۲۶	۲۲۹۰۰۰	-
۲	SEP	۰/۰۴	۰/۰۳	۸۴/۴۴	۶۶/۳۲	-
۳	MAE	۰/۰۳	۰/۰۲	۶۶/۹۱	۳۷/۷۸	-
۴	REP(%)	۵/۳۹	۵/۹۸	۸/۰۲	۶/۴۵	-
۵	RMSE	۰/۰۴	۰/۰۳	۸۴/۴۳	۶۶/۳۲	-
۶	MRE	۴/۳۴	۴/۱۴	۶/۵۳	۳/۳۲	-
۷	R <sup>2</sup>	۰/۹۲	-	۰/۸۹	-	R <sup>2</sup> > ۰/۶
۸	Q <sub>LOO</sub> <sup>2</sup>	-	۰/۹۴	-	۰/۸۸	Q <sub>LOO</sub> <sup>2</sup> > ۰/۵
۹	R <sub>0</sub> <sup>2</sup>	۰/۹۱	۰/۹۳	۰/۸۹	۰/۸۶	نزدیک به R <sup>2</sup>
۱۰	R <sub>0</sub> <sup>2</sup> نسبی	۰/۰۱	۰/۰۱	۰/۰۱	۰/۰۲	< ۰/۱
۱۱	R <sub>m</sub> <sup>2</sup>	۰/۸۳	۰/۸۵	۰/۸	۰/۷۶	> ۰/۵
۱۲	R <sub>0</sub> ' <sup>2</sup>	۰/۹۲	۰/۹۳	۰/۸۸	۰/۸۸	نزدیک به R <sup>2</sup>
۱۳	R <sub>0</sub> ' <sup>2</sup> نسبی	۰	۰/۰۱	۰	۰/۰۲	< ۰/۱
۱۴	R <sub>m</sub> ' <sup>2</sup>	۰/۸۲	۰/۹۳	۰/۷۹	۰/۷۴	> ۰/۵
۱۵	R-R	۰/۰۱	۰	۰/۰۱	۰/۰۲	< ۰/۳
۱۶	k	۰/۹۸	۱/۰۰	۱/۰۴	۰/۹۹	۰/۸۵ ≤ k ≤ ۱/۱۵
۱۷	k'	۱/۰۱	۱/۰۰	۰/۹۵	۱	۰/۸۵ ≥ k' ≤ ۱/۱۵

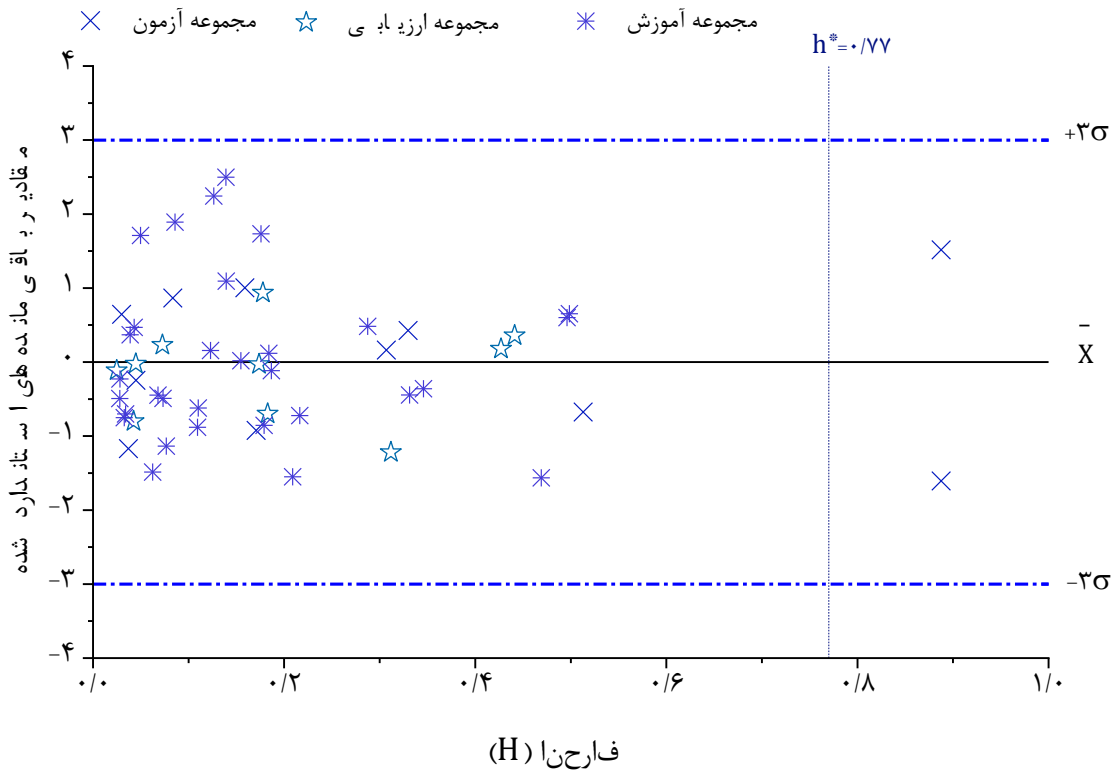


## ۲-۵-۷-۴ ارزیابی مدل SCAD-ANN با استفاده از دامنه کاربرد

هدف از دامنه کاربرد، اثبات اعتمادپذیری مدل پیشنهادی SCAD-ANN توصیه شده برای ترکیبات جدید (مجموعه آزمون)، است. بنابراین، وجود مقادیر RI پیش‌بینی‌شده در محدوده فضای شیمیایی، صحت نتایج را اثبات می‌کند. نمودار ویلیام به‌عنوان بهترین نمایش از دامنه کاربرد، مطابق با روش کار بخش‌های ۱-۵-۸-۵ و ۲-۲-۷-۴ برای مدل‌های QSRR مبتنی بر شبکه عصبی و توسعه یافته با توصیف‌کننده‌های منتخب روش SCAD به‌دست آمد. بنابراین مقادیر انحراف بر اساس رابطه ۱-۱۳ محاسبه شد. به‌این ترتیب مقادیر استاندارد شده باقی‌مانده‌ها بر حسب مقادیر پیش‌بینی شده و واقعی RI و مطابق با رابطه ۱-۱۴ محاسبه شد. از رسم مقادیر باقی‌مانده‌های استاندارد شده بر حسب مقادیر انحراف نمودار ویلیام به‌دست آمد (شکل ۲-۴۷ و شکل ۲-۴۸). علاوه بر این مقادیر  $h^*$  برای هر دو مجموعه داده A و B به‌ترتیب برابر با ۰/۴۲ و ۰/۷۷ به‌دست آمد. نمودار ویلیام نشان می‌دهد که بیش از ۹۵ درصد داده‌ها در محدوده اطمینان (کوچک‌تر از  $h^*$  و قرارگیری در محدوده  $\pm 3\sigma$ ) قرار دارند، به این معنی که نتایج پیش‌بینی شده مدل‌های SCAD-ANN هر دو مجموعه داده دارای اعتماد و استحکام قابل قبولی هستند.



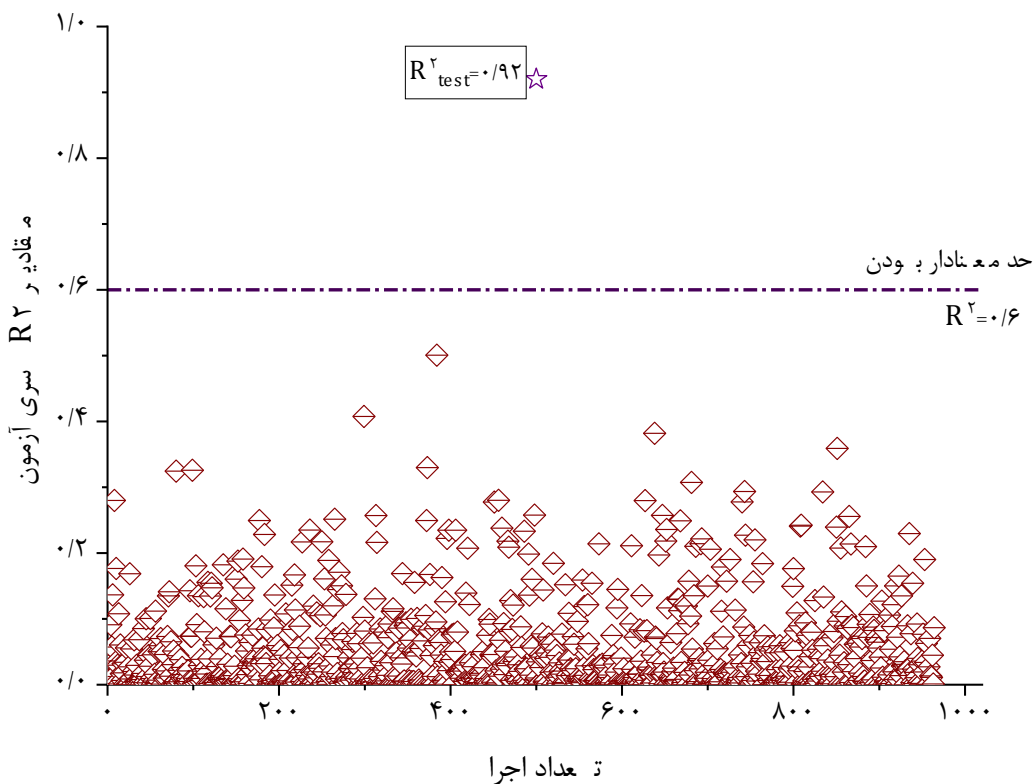
شکل ۲-۴۷ دامنه کاربرد مدل SCAD-LM-ANN برای مجموعه A، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $\pm 3\sigma$  و  $h^*$  است.



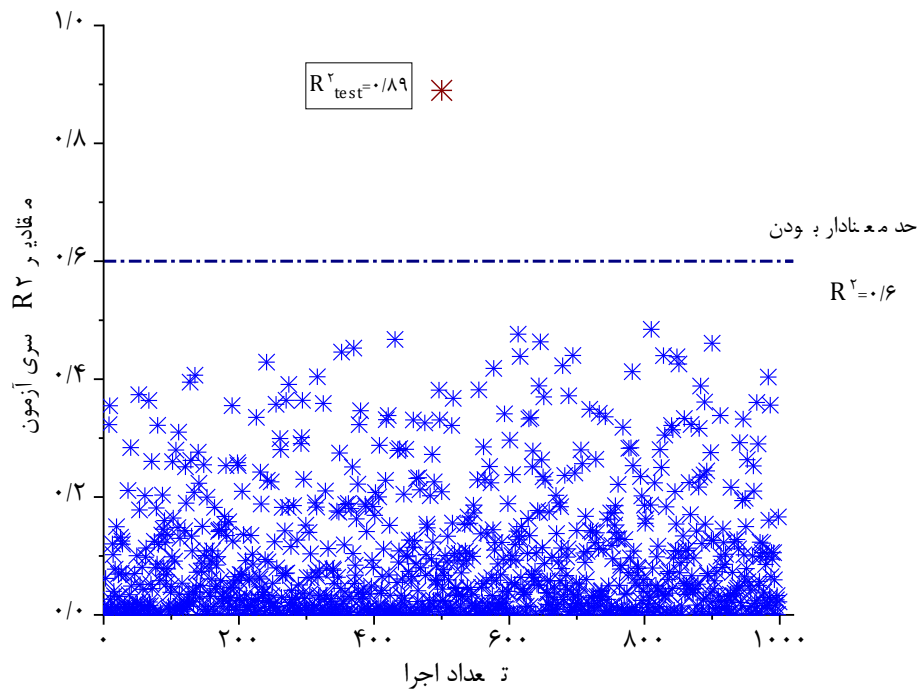
شکل ۲-۴۸ دامنه کاربرد مدل SCAD-BR-ANN برای مجموعه B، خطوط نقطه چین افقی و عمودی در دو انتهای نمودار به ترتیب نمایانگر مقادیر  $\pm 3\sigma$  و  $h^*$  است.

## ۲-۵-۷-۵ ارزیابی مدل SCAD-ANN با استفاده از آزمون Y-تصادفی

به منظور بررسی اعتبار و عدم وجود ارتباط شانس ایجاد شده توسط مدل های SCAD-ANN بین توصیف کننده های منتخب و شاخص بازداری مربوطه از آنالیز Y-تصادفی استفاده شد. برای اجرای آزمون Y-تصادفی، ابتدا مقادیر متغیر وابسته (RI) مجموعه آموزش مجموعه داده های A و B، در محدوده حداقل و حداکثر مقادیر RI، ۱۰۰۰ بار تصادفی شدند. مدل های SCAD-ANN پیشنهادی با استفاده از پاسخ های دست کاری شده، آموزش داده شد و برای پیش بینی مقادیر RI مجموعه آزمون به کار گرفته شد. نتایج  $R^2$  حاصل از پیش بینی RI ترکیبات مجموعه آزمون با ۱۰۰۰ مدل توسعه یافته با متغیر وابسته تصادفی، به دست آمد و در شکل ۲-۴۹ و شکل ۲-۵۰ نشان داده شد. نتایج به دست آمده نشان می دهد که مقادیر  $R^2$  با مدل های SCAD-ANN توسعه یافته با داده های دست کاری شده، به طور قابل توجهی کوچک تر از ضریب تعیین مجموعه آزمون مدل برتر هر دو مجموعه داده ( $R^2 = 0/92$  برای مجموعه A و  $R^2 = 0/89$  برای مجموعه B) و حتی کوچک تر از مقدار قابل قبول  $0/6$  است. بنابراین، نتایج به دست آمده ثابت می کند که رابطه QSRR ایجاد شده تصادفی نیست و بر اساس یک رابطه دقیق و منطقی بین توصیف کننده های منتخب روش SCAD و RI ایجاد شده است.



شکل ۲-۴۹ نمودار مقادیر  $R^2$  به دست آمده در آزمون  $Y$ -تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای  $Y$ -تصادفی و پیش‌بینی ی‌ب‌بی ترکیبات آزمون مجموعه  $A$  به وسیله مدل  $SCAD-ANN$  با استفاده از پاسخ تصادفی شده در شرایط بهینه



شکل ۲-۵۰ نمودار مقادیر  $R^2$  به دست آمده در آزمون  $Y$ -تصادفی بر حسب تعداد اجرا برای ۱۰۰۰ اجرای  $Y$ -تصادفی و پیش‌بینی ی‌ب‌بی ترکیبات آزمون مجموعه  $B$  به وسیله مدل  $SCAD-ANN$  با استفاده از پاسخ تصادفی شده در شرایط بهینه

# نیچہ گیری و آئندہ نگری



### ۳-۱ نتیجه گیری نهایی مدل‌های توسعه یافته QSAR/QSPR

در این رساله تلاش شد تا مدل‌های QSAR/QSPR با اعتبار قابل قبول جهت پیش‌بینی فعالیت دارویی و یا شاخص بازداري ترکیبات شیمیایی ارائه شود. ارزیابی مدل‌های توسعه یافته نشان داد که این مدل‌ها از اعتبار و تعمیم‌پذیری مناسبی برخوردار هستند. بنابراین در این فصل سعی شده است ارتباط بین توصیف‌کننده‌های مورد استفاده در مدل‌های بهینه شبکه عصبی با متغیر پاسخ هدف (فعالیت دارویی و یا شاخص بازداري) تجزیه و تحلیل گردد. به طوری که با فهم چگونگی تأثیر هر توصیف‌کننده بر متغیر پاسخ و میزان سهم مشارکت هر توصیف‌کننده در مدل برتر، ترکیبات جدیدی با فعالیت دارویی مناسب به عنوان ترکیبات کاندید برای سنتز و بررسی بیشتر به طور تجربی پیشنهاد گردد. لازم به ذکر است که برای اثبات صحت و درستی فعالیت دارویی ترکیبات پیشنهادی جدید از مطالعه داکینگ مولکولی و بررسی برهم‌کنش لیگاند - گیرنده استفاده شده است. پتانسیل بالقوه فعالیت دارویی ترکیبات پیشنهادی با مقایسه برهم‌کنش‌های مؤثر با اسیدآمینوهای کلیدی در ترکیبات فعال مجموعه داده‌ها با ترکیبات پیشنهادی به صورت نظری تأیید گردید. علاوه بر این موارد، ویژگی‌های فارماکوکینتیکی ترکیبات پیشنهادی نیز با استفاده از ابزار وب رایگان Swiss-ADME محاسبه و مورد بررسی و تحلیل قرار گرفته است. در ادامه به جزییات مطالعات ذکر شده QSAR و QSRR برای هر کدام از مجموعه داده‌ها با استفاده از مدل‌های برتر معرفی شده در فصل‌های قبل پرداخته خواهد شد.

## ۳-۱-۱ تجزیه و تحلیل توصیف‌کننده‌های مدل SCAD-LM-ANN برای مجموعه

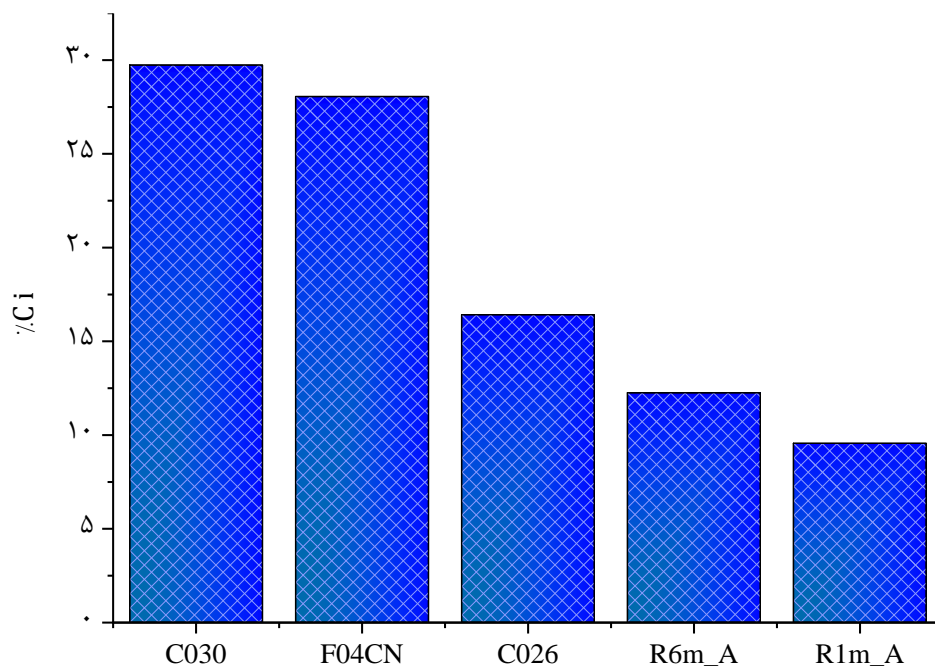
### بازدارنده‌های ایدز

#### ۳-۱-۱-۱ محاسبه سهم مشارکت هر توصیف‌کننده در مدل SCAD-ANN

به‌منظور بررسی اهمیت و میزان مشارکت هر توصیف‌کننده در مدل SCAD-LM-ANN درصد سهم هر توصیف‌کننده ( $C_i\%$ ) در مدل SCAD-ANN در شرایط بهینه، برآورد شد. روش کار به این صورت بود که برای محاسبه سهم مشارکت توصیف‌کننده  $i$ ، مقادیر آن توصیف‌کننده به‌طور تصادفی در محدوده مقادیر واقعی آن‌ها تغییر داده شد. پارامترهای شبکه عصبی (گره، دور آموزش، توابع آموزش و انتقال) مدل شبکه عصبی با استفاده از یک زیر مجموعه داده با ۵ توصیف‌کننده به‌عنوان ورودی، در مقادیر بهینه (گزارش شده در بخش ۲-۲-۶) قرار گرفتند و این مدل شبکه عصبی برای محاسبه سهم مشارکت هر توصیف‌کننده مورد استفاده قرار گرفت. مدل شبکه عصبی در شرایط بهینه هر بار با استفاده از یک زیر مجموعه با ۵ توصیف‌کننده، آموزش داده شد، در حالی که مقادیر توصیف‌کننده  $i$  با مقادیر تصادفی جایگزین شده بودند. پس از آموزش شبکه و پیش‌بینی فعالیت دارویی ترکیبات مجموعه ارزیابی مقدار MAE مجموعه ارزیابی در حالی که مقادیر توصیف‌کننده  $i$  تصادفی شده بودند، محاسبه شد ( $MAE_i$ ). این فرآیند برای هر ۵ دسته توصیف‌کننده تکرار شد و ۵ مقدار ( $MAE_i$ ) برای مجموعه ارزیابی به دست آمد. سپس، مقادیر  $C_i\%$  ( $i=1,2,\dots,5$ ) برای همه توصیف‌کننده‌ها با استفاده از رابطه ۱-۱۵ به دست آمد.

نتایج به‌دست‌آمده برای سهم مشارکت هر توصیف‌کننده در مدل SCAD-ANN در شکل ۱-۳

آورده شده است. نتایج نشان می‌دهد که سه توصیف‌کننده F04CN، C030 و C026 درصد مشارکت بالاتری در مدل بهینه SCAD-ANN دارند. بنابراین چگونگی اثر این توصیف‌کننده‌ها بر فعالیت دارویی ترکیبات مورد بررسی بیشتر قرار گرفت، تا از نتایج آن برای پیشنهاد ترکیبات جدید با فعالیت دارویی مناسب استفاده گردد.



نام تو صیف کننده های منتخب روش SCAD

شکل ۱-۳ نمودار سهم مشارکت توصیف کننده‌ها در مدل SCAD-LM-ANN

### ۳-۱-۱-۲ بررسی رابطه بین توصیف کننده‌های استفاده شده در مدل نهایی (SCAD-LM-)

#### (ANN) و فعالیت دارویی ترکیبات مورد مطالعه

بررسی میزان اثر توصیف کننده‌ها بر پاسخ (فعالیت دارویی) با مقایسه ضرایب استاندارد شده مدل خطی SCAD مورد بررسی قرار گرفت. به این منظور ضرایب استاندارد شده توصیف کننده‌های منتخب برای مدل SCAD مطابق با رابطه ۱-۳ استخراج شدند. علامت مثبت ضرایب رگرسیون (رابطه ۱-۳) نشان می‌دهد که همه توصیف کننده‌ها اثر افزایشی بر فعالیت دارویی دارند، به این معنی که افزایش مقدار توصیف کننده مورد مطالعه باعث افزایش نسبی فعالیت دارویی ترکیب می‌شود.

$$pEC_{50} = 0.95 + 0.38F04CN + 0.09 C-030 + 0.06C-026 + 0.49 R6m\_A + 0.54 R1m\_A \quad \text{رابطه ۱-۳}$$



با توجه به نمودار سهم مشارکت (شکل ۳-۱)، توصیف‌کننده‌های F04CN، C030 و C026 که بیش‌ترین مشارکت را در مدل ANN دارند و ضرایب این توصیف‌کننده‌ها در مدل SCAD (رابطه ۳-۱) مثبت است. بنابراین با افزایش مقادیر همه این توصیف‌کننده‌ها، فعالیت دارویی ( $pEC_{50}$ ) افزایش می‌یابد. F04CN به طبقه توصیف‌کننده‌های اثر انگشت فرکانس دو بعدی<sup>۱</sup> تعلق دارد و به تعداد C-N در فاصله هندسی  $4\text{\AA}$  مربوط است و تعداد کربن موجود در فاصله هندسی  $4\text{\AA}$  از اتم نیتروژن را نشان می‌دهد. بررسی عمیق‌تر ساختار ترکیبات نشان می‌دهد که در ساختارهایی با گروه‌های عاملی تقریباً یکسان، با وجود یک حلقه‌تری آزین فعالیت دارویی افزایش یابد و اغلب باعث می‌شود که مقدار توصیف‌کننده F04CN به‌طور قابل توجهی افزایش یابد. بنابراین می‌توان گفت، F04CN به حضور حلقه‌تری آزین در ساختار ترکیبات مورد مطالعه وابسته است و حضور آن در مدل QSAR با تأثیر مثبت، تأثیر وجود حلقه‌تری آزین به فعالیت دارویی را وارد مدل نهایی می‌کند.

دو توصیف‌کننده دیگری که بیش‌ترین سهم مشارکت را در مدل ANN توسعه یافته دارند، توصیف‌کننده‌های C-030 و C-026 با اثر مثبت بر فعالیت دارویی هستند. این دو توصیف‌کننده، متعلق به دسته اجزای اتم محور هستند که در ادامه بیش‌تر به توضیح آن‌ها پرداخته می‌شود.

C-030 با نماد (X-CH-X) نشان‌دهنده اتم کربن آروماتیک متصل به گروه الکترون‌گاتیو ( $X = F$ ، O، N، S، P و هالوژن‌ها) است. اگر ساختار دارای یک گروه الکترون‌گاتیو متصل به اتم کربن آروماتیک باشد، مقدار C-030 برابر با ۱ و در غیر این صورت ۰ است. در واقع، C-030 نشان‌دهنده وجود یا عدم وجود یک گروه الکترون‌گاتیو متصل به کربن با هیبریداسیون  $sp^2$  است و تأثیر مثبت حضور گروه‌های الکترون‌گاتیو در ساختار ترکیبات مورد مطالعه بر فعالیت دارویی را در مدل وارد می‌کند.

<sup>۱</sup>Frequency binary fingerprints descriptors

<sup>۲</sup>Atom-centered fragments

توصیف کننده C-026 به معنی فراوانی R-CX-R در ترکیبات مورد مطالعه است. R و X به ترتیب نشان دهنده گروه آلکیل و گروه الکترونگاتیو هستند. C-026 مربوط به تعداد هترواتم‌های الکترونگاتیو متصل به دو گروه آلکیل است. مقدار C-026 با افزایش تعداد اتم‌های الکترونگاتیو به گروه‌های آلکیل افزایش می‌یابد. ظهور C-026 در مدل نهایی با اثر مثبت نشان می‌دهد که فعالیت دارویی با افزایش تعداد گروه‌های الکترونگاتیو افزایش می‌یابد.

### ۳-۱-۱-۳ پیشنهاد ترکیبات جدید با فعالیت دارویی مناسب با استفاده از مدل SCAD-LM-ANN ارائه شده

با استفاده از توصیف‌کننده‌های موجود در مدل برتر SCAD-LM-ANN و با توجه به رابطه ۱-۳، می‌توان دریافت که تغییر در مقادیر توصیف‌کننده‌های مهم مدل QSAR می‌تواند فعالیت دارویی را افزایش دهد. بنابراین با دانستن رابطه بین مقادیر توصیف‌کننده‌ها و جزئیات ساختاری ترکیبات مورد مطالعه می‌توان ترکیبات جدید با فعالیت دارویی قابل قبول پیش بینی کرد پیدا کرد. با توجه به توضیحات بخش ۳-۱-۱-۲ مشاهده شد که فرکانس کربن در فاصله هندسی  $4\text{\AA}$  از اتم نیتروژن، وجود اتم‌های کربن متصل به گروه‌های الکترونگاتیو و فراوانی هترواتم‌های متصل به گروه‌های آلکیل، فعالیت دارویی این دسته از ترکیبات مورد مطالعه را افزایش می‌دهد. در نتیجه، با استفاده از این رویکرد، تغییرات هدفمندی در ویژگی‌های ساختاری ترکیبات مورد مطالعه ایجاد شد و ترکیبات جدیدی با فعالیت دارویی مناسب ایجاد گردید و به‌عنوان ترکیباتی با فعالیت دارویی بالقوه پیشنهاد شدند. ساختار ترکیبات پیشنهادی با استفاده از نرم‌افزار هایپرکم رسم و بهینه‌سازی شد و توصیف‌کننده‌های مدل با استفاده از نرم‌افزار دراگون محاسبه شدند. ۵ توصیف‌کننده مهم محاسبه‌شده برای ترکیبات پیشنهادی به صورت یک زیر مجموعه با ۵ توصیف‌کننده به‌عنوان ورودی در شرایط بهینه وارد مدل بهینه SCAD-LM-ANN شدند و فعالیت‌های دارویی معادل آن‌ها ( $PEC_{50}$ ) پیش‌بینی شد. ساختار ترکیبات پیشنهادی و مقادیر  $PEC_{50}$  پیش‌بینی‌شده مربوط به

آن‌ها در جدول ۳-۱ خلاصه شده‌اند. نتایج نشان می‌دهد که مقادیر  $pEC_{50}$  پیش‌بینی‌شده در محدوده فعالیت ترکیب فعال موجود در مجموعه داده‌ها (ترکیب ۴۴ با فعالیت دارویی برابر با ۷/۷۴) می‌باشد. اگرچه مقادیر  $pEC_{50}$  پیش‌بینی‌شده برای ترکیبات پیشنهادی به اندازه‌ای بزرگ هستند که به‌عنوان ترکیبات با فعالیت بالا در نظر گرفته شوند، اما ممکن است عملاً با مولکول‌های هدف بر هم کنش مناسبی نداشته باشند و فعالیت‌های مورد انتظار را نشان ندهند. بنابراین درستی پیشنهاد این ترکیبات با به‌کارگیری یک روش مناسب اثبات شد. عملاً بهترین روش برای اثبات فعالیت بالای ترکیبات پیشنهادی، سنتز آن‌ها و تعیین فعالیت تجربی فعالیت دارویی آن‌هاست. باین‌حال، این فرآیند زمان‌بر و پرهزینه است و به امکانات آزمایشگاهی نیازمند است که انجام آن در مطالعه حاضر میسر نبوده، بنابراین یک روش نظری برای تأیید فعالیت ترکیبات پیشنهادی از طریق، مقایسه برهم‌کنش‌های گیرنده-لیگاند ترکیبات پیشنهادی با ترکیبات فعال مجموعه داده‌های مورد مطالعه با استفاده از نرم‌افزار داکینگ مولکولی مورد استفاده قرار گرفت. مراحل انجام داکینگ مولکولی، استخراج برهم‌کنش‌های گیرنده-لیگاند و مقایسه این برهم‌کنش‌ها برای ترکیبات مختلف در ادامه آورده شده است.

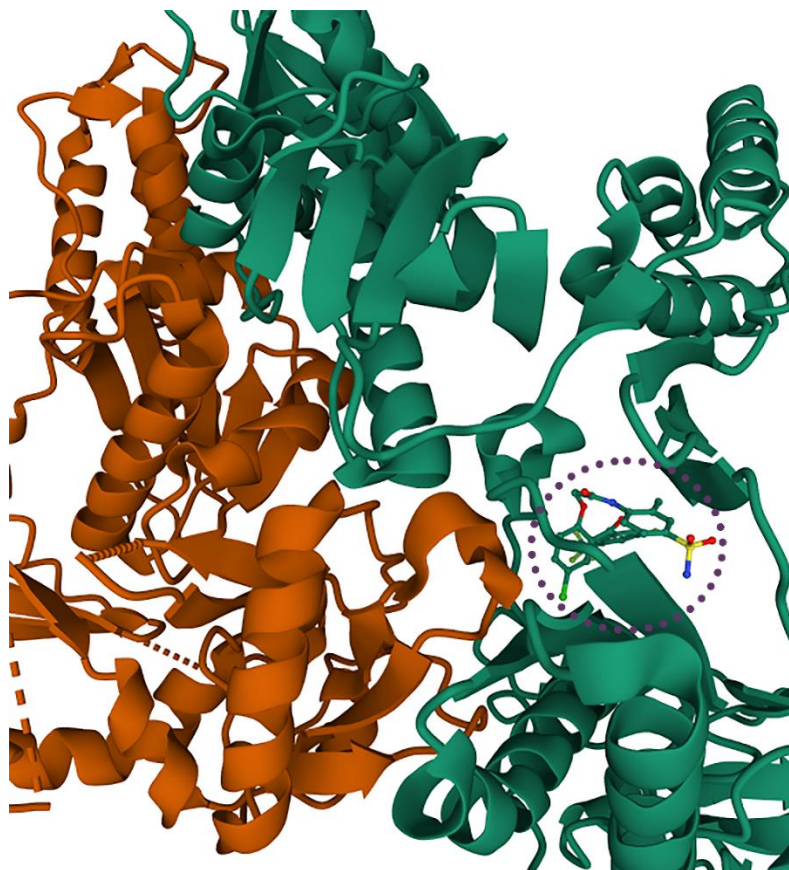
### ۳-۱-۱-۴ مطالعه داکینگ مولکولی

به‌منظور انجام عملیات داکینگ مولکولی باید مراحل متفاوتی را انجام داد که در فصل اول (بخش‌های ۱-۷ تا ۱-۷-۵) به آن‌ها اشاره شده است. ابتدا ساختار کریستالوگرافی گیرنده مورد استفاده برای ترکیبات مورد مطالعه به پیشنهاد مقالات شناسایی شد. گیرنده با کد 3DLG شناسایی و از سایت بانک اطلاعاتی پروتئین با پسوند pdb دانلود شد [۱۹۶]. ارزش تفکیک<sup>۱</sup> برای ساختار کریستالوگرافی گیرنده 3DLG برابر با  $2/20^{\circ}A$  است که برای انجام مطالعات داکینگ مولکولی مناسب است [۱۹۷]. ساختار کریستالوگرافی گیرنده 3DLG و لیگاند کریستالوگرافی موجود در زنجیره A در شکل ۳-۲ نشان داده شده

---

<sup>۱</sup>Resolution

است. در مرحله بعد ساختار کریستالوگرافی گیرنده 3DLG در نرم‌افزار ویورلایت ورژن ۵ فراخوانی شد و عملیات آماده‌سازی متفاوتی از جمله حذف مولکول‌های آب، حذف کوفاکتورها و حذف زنجیره اسید آمینه‌ای فاقد لیگاند کریستالوگرافی انجام شد و باقی‌مانده ساختار که شامل زنجیره اسید آمینه‌ای A و لیگاند کریستالوگرافی است که با پسوند pdb ذخیره شد. سپس زنجیره اسید آمینه‌ای به‌عنوان فایل گیرنده و ترکیب موجود در ساختار کریستالوگرافی به‌عنوان لیگاند به‌طور مجزا با پسوند pdb ذخیره شدند و به‌عنوان ورودی‌های نرم‌افزار داکینگ برای انجام مرحله اعتبار سنجی داکینگ مورد استفاده قرار گرفتند. در مرحله بعد، مختصات جایگاه فعال گیرنده با توجه به مختصات مرکز ثقل لیگاند کریستالوگرافی و با استفاده از نرم‌افزار ویورلایت با مختصات  $x = -2/220$ ,  $y = -34/568$ ,  $z = 22/675$  استخراج گردید.



شکل ۲-۳ ساختار کریستالوگرافی 3DLG [۱۹۸] (منطقه نقطه چین نشان‌دهنده لیگاند کریستالوگرافی و مابقی زنجیره‌های اسید آمینه‌ای است)

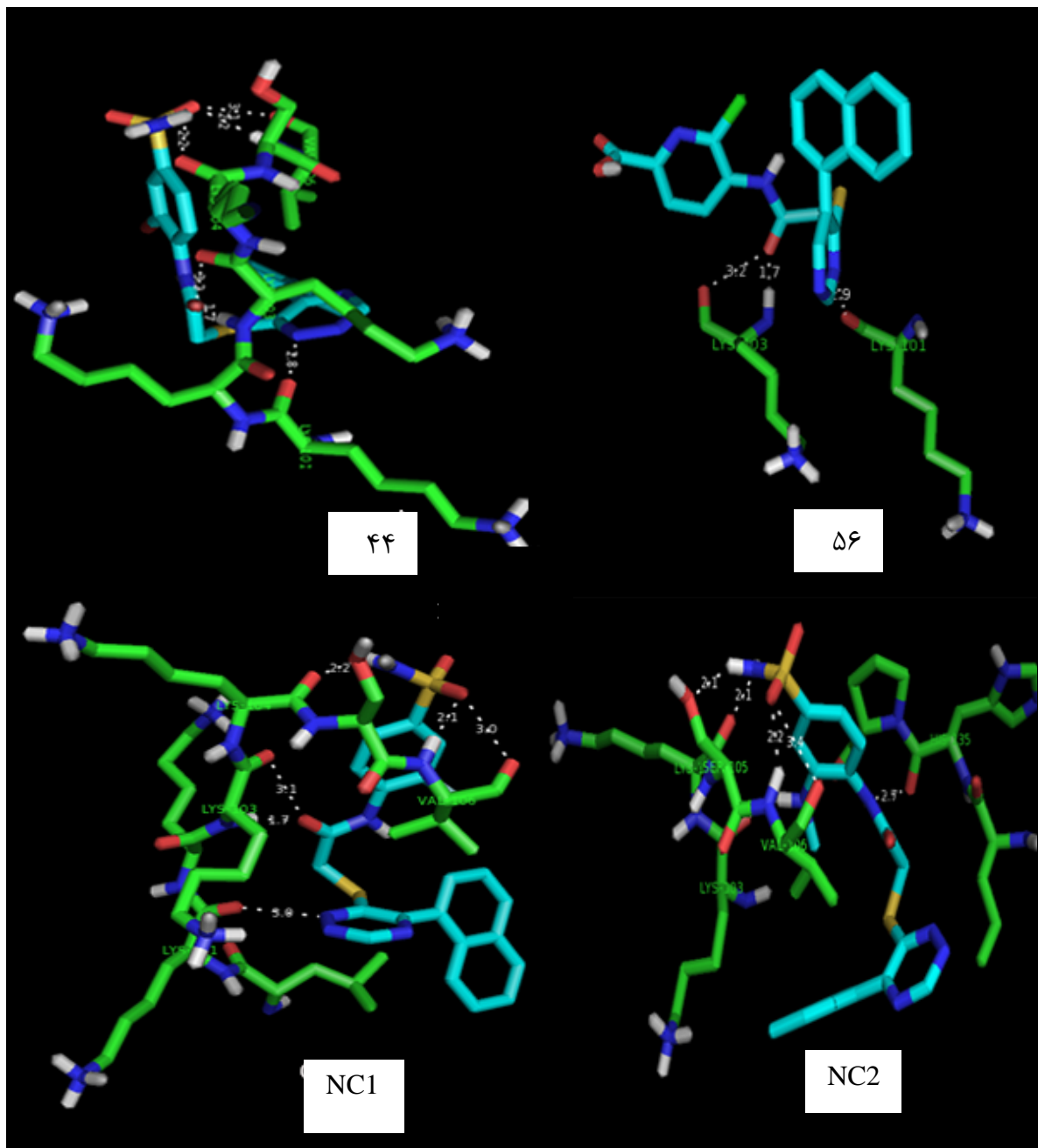
پس از آماده‌سازی فایل‌های ورودی (زنجیره اسید آمینه حاوی ساختار لیگاند کریستالوگرافی به‌عنوان گیرنده و لیگاند کریستالوگرافی به‌عنوان لیگاند) فرایند اعتبار سنجی داکینگ مولکولی با استفاده از نرم‌افزار Autodock 4.2 انجام شد. به‌این منظور ابتدا فایل گیرنده در محیط نرم‌افزار وارد شد. آماده‌سازی‌های بیش‌تر روی ساختار گیرنده از جمله افزودن هیدروژن برای جبران کمبود هیدروژن ساختار کریستالوگرافی، ادغام هیدروژن‌های غیر قطبی متصل به هر اتم کربن و افزایش بار الکتریکی کلمن انجام شد [۴۴، ۵۰، ۵۱]. در مرحله بعد به‌منظور انجام فرایند اعتبار سنجی، ساختار لیگاند کریستالوگرافی به‌عنوان لیگاند در نرم‌افزار Autodock 4.2 فراخوانی شد و با پسوند pdbqt ذخیره شد. در مرحله بعد خروجی pdbqt ساختار گیرنده نیز ذخیره شد و مختصات یک شبکه در بر گیرنده جایگاه فعال گیرنده با پیشنهاد مقالات منتشر شده در حوزه مطالعات داکینگ مولکولی، با ابعاد  $60 \times 60 \times 60 \text{ \AA}$  و فاصله بین نقاط  $0.375 \text{ \AA}$  (در حدود یک چهارم طول پیوند یگانه کربن-کربن) ایجاد شد و بر هم کنش‌های گیرنده-لیگاند در این شبکه مورد بررسی قرار گرفت [۱۹۹]. داکینگ مولکولی با استفاده از اجرای الگوریتم ژنتیک لامارکین (LGA)، در تعداد اجراهای متفاوت ۱۰۰، ۱۵۰ و ۲۰۰ انجام شد. لازم به ذکر است که سایر پارامترهای الگوریتم ژنتیک در حالات پیش‌فرض پیشنهادی نرم‌افزار در نظر گرفته شدند. عملیات اعتبارسنجی در شرایط ذکر شده، انجام شد و فایل‌های خروجی با پسوند dlg استخراج و ذخیره شدند. با توجه به نتایج داکینگ مولکولی در مرحله اعتبار سنجی در ۳ حالت با تعداد اجراهای متفاوت الگوریتم LGA، تعداد اجرای الگوریتم LGA برابر با ۱۵۰ کم‌ترین مقدار ریشه میانگین مربعات انحراف (RMSD) و بیش‌ترین تعداد کنفورماسیون در خوشه اول را دارا بوده و در نتیجه، برای انجام فرایند داکینگ مولکولی ترکیبات مورد نظر از تعداد اجرای الگوریتم LGA برابر با ۱۵۰ استفاده شد. در فرایند انجام داکینگ مولکولی ترکیبات مورد مطالعه، همانند فرایند اعتبار سنجی، تمامی شرایط (مختصات جایگاه فعال، ابعاد شبکه و پارامترهای پیش‌فرض الگوریتم

---

<sup>1</sup>Kollman

ژنتیک) در شرایط داکینگ مولکولی ترکیبات جدید نیز تعریف شد. بنابراین داکینگ مولکولی در شرایط بهینه برای ترکیبات ۴۴ (با بیشترین فعالیت ( $pEC_{50} = 7/74$ ), ترکیب ۵۶ با کمترین فعالیت ( $4/62$ )  $pEC_{50} =$ ) و تمامی ترکیبات پیشنهادی در جایگاه فعال گیرنده انجام شدند. بهترین پیکربندی‌ها از ترکیبات مورد مطالعه و پیشنهادی با توجه به حداقل انرژی آزاد اتصال و بیشترین تعداد پیکربندی در خوشه اول، شناسایی شدند و برهم‌کنش آن‌ها با اسید آمینه‌های کلیدی با استفاده از نرم‌افزار vmd مورد بررسی قرار گرفت. مطابق با شکل ۳-۳ و شکل ۴-۳ و داده‌های جدول ۱-۳، فعال‌ترین ترکیب (ترکیب شماره ۴۴) دارای برهم‌کنش‌های آب‌دوست با اسیدهای آمینه کلیدی Val106 (2.2 Å), Val106 (3.1 Å), Lys104 (2.2 Å), Lys103 (3.2 Å), Lys103 (1.7 Å) و Lys101 (2.8 Å) می‌باشد. درحالی‌که ضعیف‌ترین ترکیب (ترکیب شماره ۵۶) دارای تعداد برهم‌کنش‌های کم‌تر با تعداد محدودتری اسید آمینه Lys101 (2.9 Å), Lys103 (3.2 Å) و Lys103 (1.7 Å) است (موارد داخل پرانتز طول پیوند اسید آمینه کلیدی با یک اتم در ساختار مورد مطالعه را نشان می‌دهد). بنابراین ترکیب فعال (ترکیب شماره ۴۴) پیوندهای هیدروژنی مناسبی را با اسید آمینه‌ها برقرار کرده است، که به‌طور معناداری از نظر برهم‌کنش متفاوت از ترکیب ضعیف‌تر (ترکیب شماره ۵۶) است. این به آن معناست که برهم‌کنش‌های حاصل از داکینگ مولکولی معیار مناسبی برای ارزیابی میزان فعالیت دارویی ترکیبات مورد مطالعه هستند. بنابراین برهم‌کنش‌های لیگاند-گیرنده برای تمام ترکیبات پیشنهادی استخراج و مورد تجزیه و تحلیل قرار گرفت و نتایج در شکل ۳-۳ و شکل ۴-۳ آورده شده است. نتایج به‌وضوح نشان می‌دهد که ترکیبات پیشنهادی پیوندهای هیدروژنی مناسبی را با اسید آمینه‌های کلیدی برقرار نموده‌اند. در نتیجه فعالیت دارویی پیش‌بینی شده برای آن‌ها به‌وسیله مدل به‌خوبی برآورد شده‌اند. علاوه بر مطالعه داکینگ مولکولی، پارامترهای قاعده لیپینسکی نیز برای ترکیبات پیشنهادی محاسبه شدند و نتایج در جدول ۱-۳ خلاصه شده‌اند. نتایج نشان می‌دهد که تمامی پارامترهای مورد نظر از جمله وزن مولکولی ( $MW < 500$ ), چربی‌دوستی ( $MLOGP < 4/15$ ), تعداد

گیرنده‌های پیوند هیدروژنی ( $\#H-B-acc < 10$ )، تعداد دهنده‌های پیوند هیدروژنی ( $\#H-B-don < 5$ ) و تعداد پیوندهای قابل چرخش ( $\#Rot-B < 10$ ) برای ترکیبات پیشنهادی دارای مقادیر قابل قبولی هستند [۲۰۰]. بنابراین ترکیبات پیشنهادی معیارهای شباهت به ترکیبات دارویی را داشته‌اند و به‌عنوان ترکیبات بالقوه فعالی که شباهت به ترکیب رهبر دارند، می‌توانند به‌عنوان کاندیدهای فعال مناسبی برای سنتز و آزمایش‌های بیولوژیکی مورد توجه قرار گیرند. علاوه بر این فاکتورها، پارامتر سهولت سنتز (Syn-Acc) نیز با استفاده از سایت Swiss-ADME [۲۰۰] استخراج شد و در جدول ۳-۱ آورده شد. فاکتور سهولت سنتز می‌تواند دارای مقداری در محدوده ۱ (سنتز بسیار آسان) تا ۱۰ (ترکیب پیچیده و چالش‌برانگیز) باشد [۲۰۱]. فاکتور سهولت سنتز برای ترکیبات پیشنهادی دارای مقدار قابل قبولی می‌باشد و در نتیجه سنتز آن‌ها در محیط آزمایشگاهی امکان‌پذیر خواهد بود.

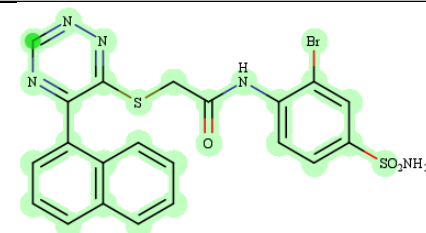
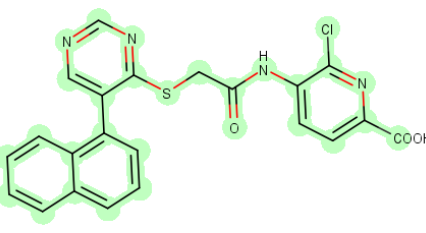
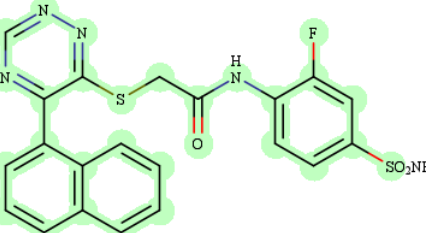
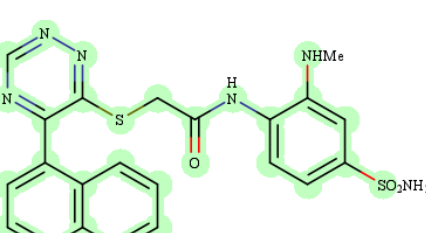
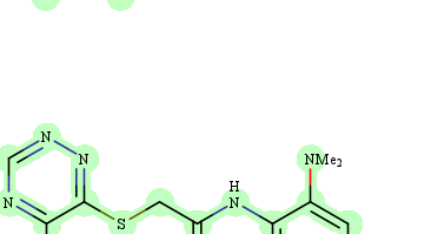


شکل ۳-۳ بررسی برهم‌کنش ترکیبات ۴۴، ۵۶، NC1 و NC2 با گیرنده

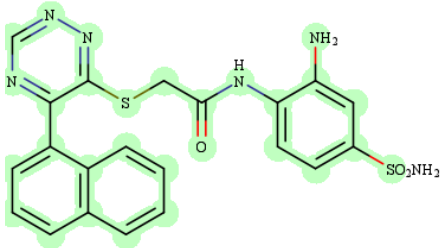
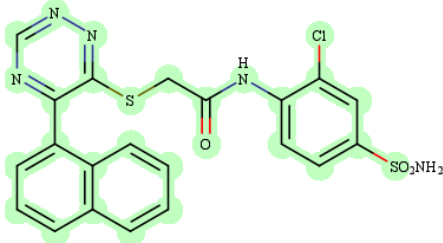
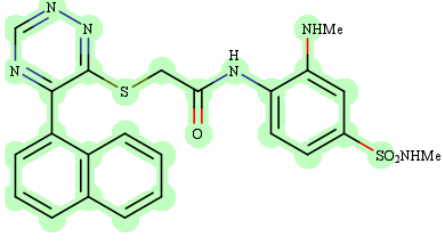
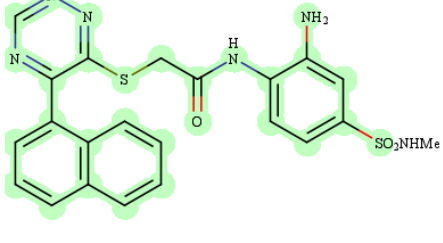
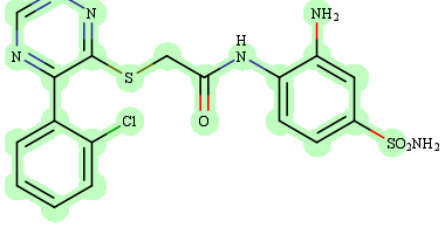




جدول ۱-۳ پارامترهای PK محاسبه شده برای ترکیبات مورد مطالعه و ترکیبات پیشنهادی

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
۴۴		۵۳۰/۴	۱/۹۴	۷	۷	۲	۳/۳۹	۷/۷۴
۵۶		۴۵۰/۹	۰/۹۸	۷	۶	۲	۳/۲۵	۴/۶۲
NC1		۴۶۶/۵	۰/۸۴	۷	۷	۳	۳/۴۸	۷/۸۲
NC2		۴۸۶/۰	۱/۸۳	۷	۷	۲	۳/۴۲	۷/۸۱
NC3		۴۸۰/۶	۱/۰۵	۸	۷	۳	۳/۶۶	۷/۸۱

ادامه جدول ۱-۳

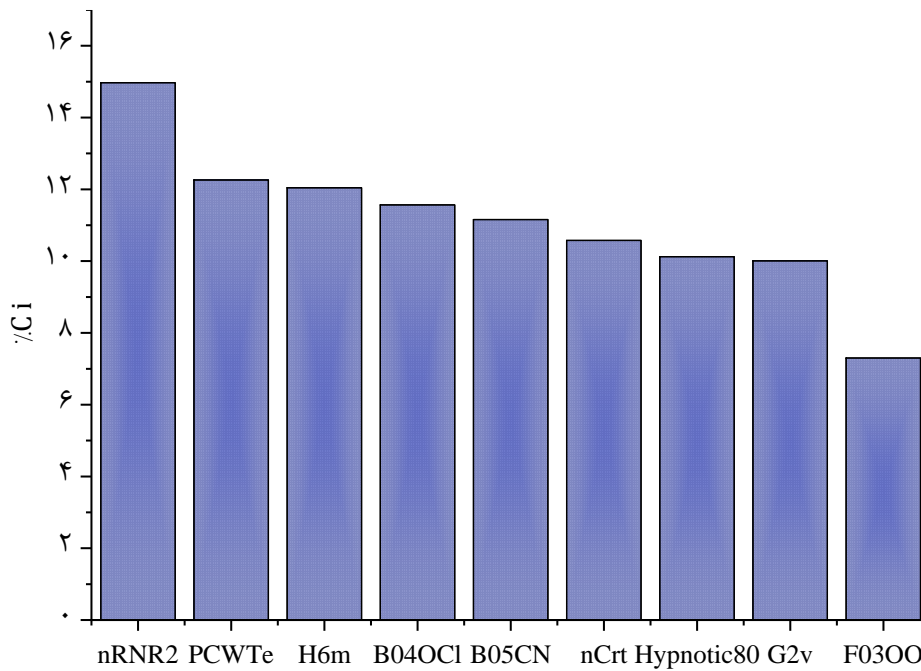
شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
NC4		۴۹۴/۶	۱/۲۷	۹	۷	۳	۳/۷۷	۷/۸
NC5		۴۸۰/۶	۱/۰۵	۸	۷	۳	۳/۵۹	۷/۸
NC6		۴۶۹/۵	۱/۷۲	۷	۸	۲	۳/۳۷	۷/۸
NC7		۴۹۴/۶	۱/۲۷	۸	۷	۲	۳/۷۷	۷/۷۹
NC8		۴۵۰/۹	۰/۵۹	۷	۷	۳	۳/۳۲	۷/۷۸

## ۳-۱-۲ تجزیه و تحلیل توصیف کننده‌های مدل ALASSO-ANN برای بازدارنده‌های

### SARS-COV-2

#### ۳-۱-۲-۱ محاسبه سهم مشارکت هر توصیف کننده در مدل ALASSO-ANN

مدل ALASSO-LM-ANN ساخته شده با استفاده از ۹ توصیف کننده مؤثر شامل (B04[O-CI])، H6m، nRNR2، PCWTe، Hypnotic80، B05[C-N]، F03[O-O]، nCrt، G2v به‌عنوان مدل برتر در نظر گرفته شد. سهم هر توصیف کننده با محاسبه درصد مشارکت توصیف کننده  $i$  ( $C_i$ ) در مدل بهینه ALASSO-LM-ANN با معماری ۱-۲-۹ در شرایط بهینه برآورد شد. به این منظور مقادیر توصیف کننده  $i$  در محدوده مقادیر واقعی به صورت تصادفی تغییر داده شد و مقادیر بقیه توصیف کننده‌ها دارای مقادیر اصلی خود بودند. مدل ALASSO-LM-ANN در شرایط بهینه با استفاده از زیر مجموعه داده‌های ایجاد شده شامل توصیف کننده  $i$  با مقادیر تصادفی و ۸ توصیف کننده دیگر با مقادیر واقعی خود آموزش داده شد و برای پیش بینی فعالیت دارویی مجموعه ارزیابی به کار گرفته شد. با استفاده از نتایج حاصله، مشابه همین فرایند با تصادفی‌سازی تمامی توصیف کننده‌ها تکرار شد. در نهایت ۹ مقدار  $MAE_i$  به دست آمد. درصد سهم مشارکت هر توصیف کننده ( $C_i$ ) طبق رابطه ۱-۱۵ محاسبه شد. نتایج به دست آمده برای درصد مشارکت توصیف کننده‌ها در مدل ALASSO-LM-ANN در شکل ۳-۵ آورده شده است. نتایج نشان می‌دهد که، مدل ALASSO-LM-ANN در حضور توصیف کننده‌هایی که دارای مقادیر  $C_i$  بیش‌تر می‌باشند، خطای پیش بینی بیش‌تری را در مدل بهینه ایجاد می‌کند. بنابراین این توصیف کننده بیش‌ترین مشارکت را در مدل ALASSO-LM-ANN به‌عنوان مدل QSAR برتر پیشنهادی دارند. در ادامه به بررسی میزان چگونگی تأثیر این توصیف کننده‌ها بر فعالیت دارویی ترکیبات و رابطه آن‌ها با ساختار ترکیبات مورد مطالعه پرداخته خواهد شد.



شکل ۳-۵ نمودار سهم مشارکت توصیف کننده‌ها در مدل ALASSO-LM-ANN  
 نام تو صیف کننده های منتخب روش ALASSO

### ۳-۱-۲ بررسی رابطه بین توصیف کننده‌های استفاده شده در مدل نهایی (ALASSO-LM-)

#### (ANN) و فعالیت دارویی ترکیبات مورد مطالعه

برای بررسی بیش‌تر چگونگی تأثیر هر توصیف کننده بر فعالیت دارویی، علامت تأثیر هر توصیف کننده بر فعالیت دارویی تعیین شد. برای این کار، ضرایب استاندارد شده مدل ALASSO با رگرسیون پاسخ (pIC<sub>50</sub>) بر حسب مقادیر توصیف کننده‌های منتخب مطابق با معادله زیر به دست آمد:

$$pIC_{50} = 0.49 + 0.66 B04[O-Cl] - 0.12 nRNR2 + 0.11 H6m - 0.09 PCWTe + 0.08 Hypnotic80$$

$$+ 0.07 B05[C-N] + 0.04 F03[O-O] + 0.03 nCrT + 0.02 G2V$$

رابطه ۳-۲

با توجه به رابطه ۳-۲، علامت ضرایب هفت توصیف کننده (Hypnotic80, H6m, B04[O-Cl])، علامت ضرایب دو مورد دیگر (PCWTe و nRNR2) دارای اثر با علامت منفی بوده و تأثیر منفی بر فعالیت دارویی دارند، در حالی که ضرایب (G2v و nCrT, F03[OO], B05[CN]) مثبت بوده و اثر افزایشی بر فعالیت دارویی دارند.

دارند. با توجه به این که توصیف کننده‌هایی مانند B04[O-Cl], mRNR2, F03[OO] و B05[CN] تأثیر بیش‌تری بر فعالیت دارویی دارند (شکل ۳-۵)، و به‌علاوه رابطه آن‌ها با ساختار شیمیایی ترکیبات تفسیر پذیرتر می‌باشند، این توصیف کننده‌ها برای اصلاح ساختارهای ترکیبات ضعیف به‌منظور یافتن ترکیباتی با فعالیت دارویی بیش‌تر (قوی‌تر) استفاده شدند و مهارکننده‌های جدید 3CL<sup>PRO</sup> با فعالیت دارویی (pIC<sub>50</sub>) مناسب پیشنهاد داده شدند. در ادامه شرح مفصل‌تری از این توصیف کننده‌ها و چگونگی وابستگی فعالیت دارویی به آن‌ها ارائه شده است.

توصیف کننده B04[O-Cl] به طبقه توصیف کننده‌های اثرات انگشت دو بعدی<sup>۱</sup> تعلق دارد و فرکانس O-Cl در فاصله مکانی ۴ Å از مرکز است. مقادیر توصیف کننده B04[O-Cl] اعداد باینری (۰ و ۱) است و در ترکیبات مورد مطالعه به وجود یا عدم وجود گروه کلروپیریدین در مجاورت گروه استر بستگی دارد. در حضور کلروپیریدین، B04[O-Cl] برابر با مقدار ۱ و در غیر این صورت برابر با صفر است. توصیف کننده B04[O-Cl] با علامت تأثیر مثبت در مدل نهایی نشان می‌دهد که وجود گروه کلروپیریدین در مجاورت گروه استر در ساختار ترکیبات سبب بهبود فعالیت دارویی می‌شود. به‌عنوان مثال، ترکیب ۷۴ دارای B04[O-Cl] برابر با ۱ (pIC<sub>50</sub> = ۷/۲۲) است، درحالی که ترکیب ۱ دارای B04[O-Cl] برابر با ۰ (pIC<sub>50</sub> = ۵/۱۴) است. بنابراین با اصلاح ساختاری و ایجاد ترکیبی با ساختاری که توصیف کننده B04[O-Cl] آن دارای مقدار ۱ باشد، ترکیب پیشنهادی فعالیت دارویی مناسبی خواهد داشت. به‌این ترتیب به‌کمک تغییر در مقدار این توصیف کننده می‌توان بازدارنده‌های جدید 3CL<sup>PRO</sup> با ویژگی‌های فارموکوکینتیک مناسب را پیشنهاد داد.

توصیف کننده nRNR2 به طبقه شمارش گروه عاملی تعلق دارد و تعداد آمین‌های آلیفاتیک نوع سوم را شمارش می‌کند. در ترکیبات مورد مطالعه، مقدار nRNR2 با حضور بخش پیرازین کنترل می‌شود.

---

<sup>1</sup>2D binary fingerprints

حضور گروه پیپرازین به دلیل اتصال اتم‌های N به گروه‌های آلیفاتیک باعث افزایش مقدار این توصیف کننده می‌شود، در حالی که فعالیت دارویی کاهش می‌یابد (به‌عنوان مثال ترکیبات ۳۳-۳۶ و ۴۴-۴۷ با  $pIC_{50}$  زیر ۵/۰۰). از طرف دیگر، nRNR2 به دلیل داشتن تأثیر با علامت منفی حضور آن در مدل باعث کاهش فعالیت دارویی می‌شود. می‌توان نتیجه‌گیری کرد که فعالیت دارویی ترکیبات پیشنهادی در غیاب گروه عاملی nRNR2، مانند گروه پیپرازین، به اندازه کافی بزرگ است و از این نتیجه‌گیری می‌توان به‌عنوان یک مفهوم اساسی در طراحی بازدارنده‌های جدید 3CL<sup>PRO</sup> استفاده کرد.

توصیف کننده F03[O-O] نیز همانند توصیف کننده B04[O-Cl] به طبقه توصیف کننده‌های اثرات انگشت دو بعدی تعلق دارد و انعکاس دهنده حضور و فرکانس گروه O-O در فاصله مکانی ۳ Å از مرکز است. با توجه به علامت مثبت ضریب تأثیر توصیف کننده F03[O-O] می‌توان گفت برای یک ترکیب هر چه مقدار فراوانی این توصیف کننده بیش‌تر باشد، فعالیت دارویی آن نیز بیش‌تر خواهد شد. با بررسی اجمالی کل مجموعه داده‌ها، می‌توان گفت که در ترکیبات با گروه فوران در مجاورت گروه کربونیل در موقعیت  $\beta$ -استر توصیف کننده F03[O-O] دارای مقدار ۲ (حداکثر مقدار آن) می‌باشند و این ترکیبات از دسته فعال‌ترین ترکیبات موجود در مجموعه داده‌ها هستند. به‌طور مشابه، ترکیبات با دی-اون روی حلقه و در مجاورت گروه کربونیل دارای F03[O-O] برابر با ۱ بوده و از نظر فعالیت دارویی فعال‌تر از ترکیباتی هستند که مقدار توصیف کننده F03[O-O] برای آن‌ها صفر است. بنابراین، این توصیف کننده تأثیر حضور گروه‌های دی-اون و حلقه‌های فوران مجاور گروه کربونیل را با علامت مثبت وارد مدل نهایی QSAR می‌کند. توصیف کننده B05[C-N] توصیف کننده مهم دیگری با تأثیر مثبت بر فعالیت دارویی است. B05[CN] مربوط به حضور یا عدم حضور کربن و نیتروژن در فاصله مکانی ۵ Å است. این توصیف کننده اهمیت وجود پیوند CN مانند سیانید، گروه‌های نیترو و آمین در فاصله مناسب در اطراف حلقه فنیل را نشان می‌دهد.

### ۳-۱-۲-۳ کاربرد مدل ALASSO-LM-ANN در طراحی و پیشنهاد ترکیبات فعال با اثر ضد

#### کووید-۱۹

یکی از کلیدی‌ترین جنبه‌های کاربردی مطالعات QSAR که اخیراً مورد توجه برخی از محققین قرار گرفته است، استفاده از مدل‌های توسعه یافته برای پیشنهاد ترکیبات جدید و در نتیجه تسهیل فرآیند طراحی ترکیبات دارویی بالقوه است. توصیف‌کننده‌های محدود موجود در مدل نهایی ALASSO-ANN می‌توانند به‌طور مؤثر نشان‌دهنده رابطه بین فعالیت دارویی ترکیبات مورد مطالعه و ویژگی‌های ساختاری آن‌ها باشد. با توجه به این‌که توصیف‌کننده‌های موجود در مدل ارائه شده در عین قابل تفسیرپذیری و ساده بودن، قدرت پیش‌بینی قابل قبولی نیز دارد، از مدل ارائه شده به‌همراه نوع رابطه بین ساختار و مقدار توصیف‌کننده‌ها برای پیشنهاد بازدارنده‌های جدید  $3CL^{pro}$  با فعالیت‌های دارویی مطلوب استفاده گردید. قابل توجه است که در بین همه توصیف‌کننده‌های مؤثر، توصیف‌کننده‌های  $F03[OO]$ ،  $nRNR2$ ،  $B04[O-Cl]$  و  $B05[CN]$  دارای تفسیر و معنای شیمیایی واضح و وابسته به ساختار بوده و همچنین رابطه فعالیت ترکیبات با این توصیف‌کننده در مدل نهایی نمایان است. بنابراین، در طراحی ترکیبات جدید این توصیف‌کننده‌ها چگونگی ارتباط اصلاحات ساختارهای مورد مطالعه و تأثیر تغییرات ایجاد شده ساختاری بر افزایش فعالیت دارویی را در اختیار قرار می‌دهند. از این‌رو برای طراحی ترکیبات جدید با فعالیت دارویی مناسب، از معنای شیمیایی این توصیف‌کننده‌ها برای اصلاح ساختاری ترکیبات در راستای ایجاد ترکیبات با فعالیت دارویی مناسب، استفاده شد. ساختار ترکیبات پیشنهادی و مقادیر  $pIC_5$  پیش‌بینی شده مرتبط به آن‌ها در جدول ۲-۳ آورده شده است. سپس، ساختارهای سه بعدی ترکیبات پیشنهادی با استفاده از برنامه هایپرکم رسم و بهینه شدند و توصیف‌کننده‌های مؤثر با استفاده از نرم‌افزار دراگون محاسبه شدند. مقادیر توصیف‌کننده‌های ترکیبات پیشنهادی (به‌عنوان یک مجموعه آزمون) در مدل ALASSO-LM-ANN در شرایط بهینه به‌عنوان ورودی تعریف شدند و مقادیر فعالیت دارویی ( $pIC_5$ ) آن‌ها پیش‌بینی شد. با توجه به مقادیر  $pIC_5$



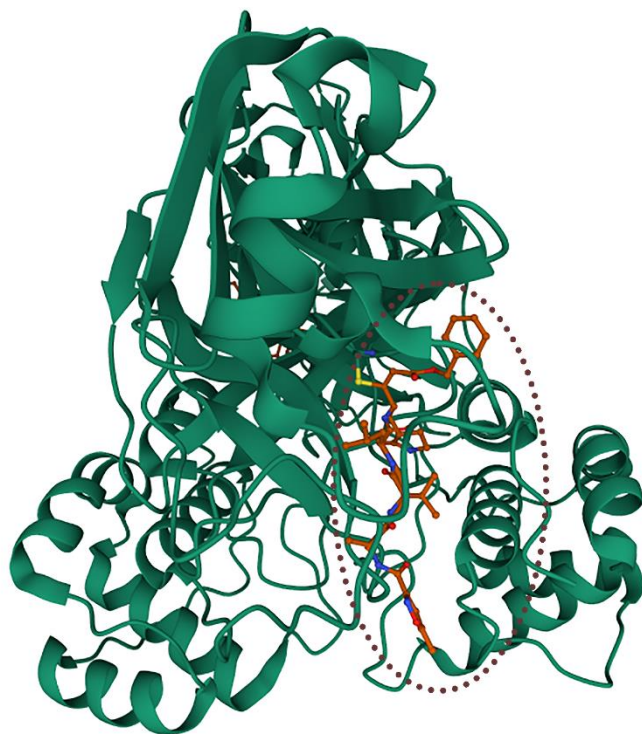
پیش‌بینی‌شده، ترکیبات پیشنهادی (NC۱ - NC۱۷) به‌عنوان بازدارنده‌های پیشنهادی بسیار فعال و ترکیبات پیشنهادی (NC۱۸ - NC۳۰) به‌عنوان بازدارنده‌های پیشنهادی نسبتاً فعال این مطالعه معرفی شدند. در ادامه برای اثبات درستی فعالیت پیش‌بینی شده برای ترکیبات فعال پیشنهادی، برهم‌کنش ترکیبات در جایگاه فعال گیرنده آنزیمی با استفاده از نرم‌افزار داکینگ مولکولی مورد بررسی قرار گرفت که در ادامه چگونگی انجام کار آورده شده است.

### ۳-۱-۲-۴ مطالعه داکینگ مولکولی بازدارنده‌های 3CL<sup>pro</sup>

صحت فعالیت دارویی ترکیبات پیشنهادی با استفاده از مطالعه داکینگ مولکولی و محاسبه خواص فارماکوکینتیک (PK) با استفاده از ابزار وب رایگان Swiss-ADME مورد بررسی قرار گرفت. در این راستا، مطالعه داکینگ همه ترکیبات پیشنهادی (NC۱ - NC۳۰) و ترکیب فعال (ترکیب ۷۴) و غیرفعال (ترکیب ۷۲) مجموعه داده‌های مورد مطالعه به‌صورت جداگانه با اسید آمینه‌های کلیدی در جایگاه فعال گیرنده پروتئینی انجام شد. برای اجرای داکینگ مولکولی ابتدا، با توجه به توصیه مقالات منتشر شده، گیرنده پروتئینی با کد کریستالوگرافی با کد 6LU7 از سایت بانک اطلاعاتی پروتئین استخراج شد [۲۰۲]. ساختار کریستالوگرافی 6LU7 با ارزش تفکیک برابر با  $2/16 \text{ \AA}$  برای انجام محاسبات داکینگ مولکولی از کیفیت مناسبی برخوردار است. شکل ۳-۶ زنجیره اسید آمینه‌ای ساختار کریستالوگرافی 6LU7 و لیگاند کریستالوگرافی موجود در زنجیره آن نشان داده شده است. اجرای داکینگ مولکولی، طبق مراحل ذکر شده در بخش‌های ۱-۷ تا ۱-۷-۵ انجام شد. قبل از داکینگ ترکیبات پیشنهادی در جایگاه فعال گیرنده، فرایند اعتبار سنجی داکینگ انجام شد. همانند قبل در انجام این فرایند ابتدا داکینگ لیگاند کریستالوگرافی در جایگاه فعال گیرنده اجرا شد و شرایط بهینه داکینگ استخراج گردید. به‌این منظور برای آماده‌سازی گیرنده، فایل داندلود شده 6LU7 با پسوند pdb در نرم‌افزار ویورلایت نسخه ۵ فراخوانی شد. مولکول‌های آب، کوفاکتورها، زنجیره‌های جانبی حذف شدند. ساختار باقی‌مانده (شامل زنجیره اسید آمینه اصلی و لیگاند

کریستالوگرافی)، به دو شکل ذخیره شد. یک بار فقط زنجیره اسید آمینه‌ای به‌عنوان ماکرومولکول ورودی نرم‌افزار Autodock4.2 و با پسوند pdb ذخیره شد و بار دیگر لیگاند کریستالوگرافی به‌طور مجزا با پسوند pdb ذخیره شد. در ادامه به‌منظور انجام فرایند اعتبارسنجی داکینگ، گیرنده در نرم‌افزار Autodock4.2 فراخوانی شد. برای آماده‌سازی بیش‌تر گیرنده، اتم‌های هیدروژن و اتم‌های هیدروژن غیر قطبی به ترتیب اضافه و ادغام شدند. سپس بار کولمن نیز برای موازنه بار سیستم اضافه شد. در مرحله بعد لیگاند کریستالوگرافی به جایگاه فعال گیرنده وارد و خروجی آن با پسوند pdbqt ذخیره شد. سپس ساختار گیرنده اصلاح شده نیز با فرمت pdbqt ذخیره شد. در مرحله بعد، یک جعبه شبکه‌ای با ابعاد  $60 \times 60 \times 60 \text{ \AA}$  و فاصله بین نقاط  $0.375 \text{ \AA}$  ایجاد شد. مرکز مختصات (x, y, z) جایگاه فعال با توجه به مختصات مرکز ثقل لیگاند برابر با  $x = -10.729 \text{ \AA}$ ,  $y = 12.418 \text{ \AA}$ ,  $z = 68.816 \text{ \AA}$ ، در نرم‌افزار Autodock4.2 تنظیم شد. سایر پارامترهای داکینگ مولکولی در مقادیر پیش‌فرض نرم‌افزار تنظیم شدند و فرایند داکینگ مولکولی با استفاده از الگوریتم ژنتیک لامارکین (LGA) در تعداد اجراهای متفاوت الگوریتم (۱۰۰، ۱۵۰، ۲۰۰) اجرا شد. با توجه به خروجی‌های داکینگ در ۳ اجرای متفاوت الگوریتم LGA، داکینگ با ۱۵۰ اجرا دارای کم‌ترین مقدار RMSD بود. در نتیجه این تعداد اجرا به‌عنوان تعداد اجرای بهینه برای اجرای فرایند داکینگ سایر ترکیبات در جایگاه فعال گیرنده ۶ULV در نظر گرفته شد. سپس به‌طور مجزا داکینگ مولکولی برای ساختار شیمیایی ترکیب شماره ۷۴ ( $pIC_{50} = 7/22$ )، به‌عنوان فعال‌ترین ترکیب، ترکیب شماره ۷۲ ( $4/00$ )  $pIC_{50} =$ ، به‌عنوان یک ترکیب با فعالیت دارویی کم و ترکیبات پیشنهادی، به‌عنوان بازدارنده‌های جدید  $3CL^{pro}$ ، در شرایط بهینه و با پارامترهای مشابه فرایند اعتبارسنجی داکینگ انجام شد. خروجی فرایند داکینگ مولکولی به‌صورت یک فایل با فرمت dlg ذخیره شد. اطلاعات مربوط به بهترین پیکربندی لیگاند موردنظر با توجه به کم‌ترین مقدار RMSD، کم‌ترین انرژی آزاد اتصال، کمترین تعداد خوشه و بیشترین تعداد کنفورماسیون در خوشه اول از خروجی dlg داکینگ مولکولی استخراج شد. برهم‌کنش مربوط به

بهترین پیکربندی در جایگاه فعال گیرنده برای همه ترکیبات مورد مطالعه در داکینگ مولکولی با استفاده از نرم افزار vmd به دست آمد. نتایج در شکل ۷-۳ تا شکل ۱۱-۳ آورده شده است.



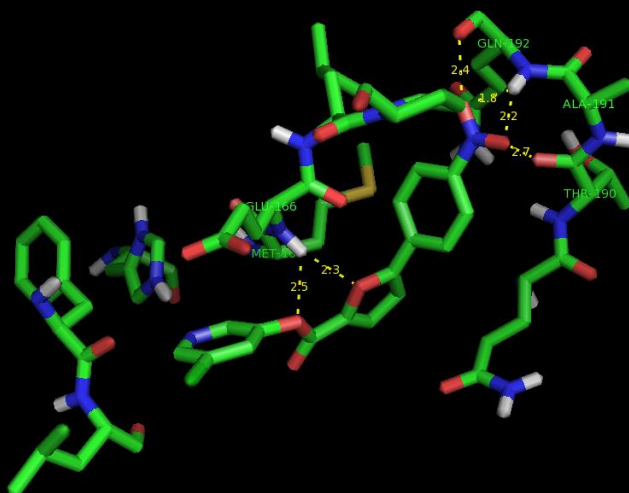
شکل ۶-۳ ساختار کریستالوگرافی 6LU7 [۱۹۸] (منطقه نقطه چین نشان دهنده لیگاند کریستالوگرافی و مابقی زنجیره های اسید آمینه ای است)

با توجه به مقالات منتشر شده، مشاهده شده است که بیش تر برهم کنش های (آب دوست، آب گریز، و اندروالسی و ...) مربوط به بازدارنده های 3CL<sup>PRO</sup> - گیرنده، با اسید آمینه های Thr24، Thr26، Phe140، Asn142، Gly143، Cys145، His163، His164، Glu166، Thr19، Thr24، His19، His163، His164، Thr19، His19، Glu166 برقرار شده است [۲۰۳] و نتایج همچنین با مشاهده متون اخیر مطابقت دارد. علاوه بر این، از بین اسید آمینه های کلیدی نام برده شده، اسید آمینه های Ala191، Asn142، Glu166، Phe140، Leu141 پیوندهای هیدروژنی مناسبی را با بازدارنده های 3CL<sup>PRO</sup> برقرار می کنند [۲۰۴، ۲۰۵]. نتایج برهم کنش ترکیبات موجود در مجموعه داده ها و ترکیبات پیشنهادی با گیرنده برای همه ترکیبات با استفاده از نتایج داکینگ مولکولی و به کمک نرم افزار PyMOL استخراج شد. نتایج برهم کنش ها در شکل

۷-۳ تا شکل ۱۱-۳ آورده شد. شکل ۷-۳ نشان می‌دهد که ترکیب ۷۴ به‌عنوان فعال‌ترین ترکیب به‌خوبی قادر به تشکیل پیوندهای هیدروژنی مناسب با Glu166، Gln192 Ala191 و Thr190 شده است. از طرفی، ترکیب ضعیف این مطالعه (ترکیب ۷۲) با کم‌ترین فعالیت دارویی تنها یک پیوند هیدروژنی با Ser144 برقرار کرده است. علاوه بر بررسی برهم‌کنش ترکیبات فعال و غیر فعال (۷۴ و ۷۲) با گیرنده، ترکیبات پیشنهادی نیز با استفاده از مطالعه داکینگ مولکولی مورد بررسی قرار گرفتند. نتایج برهم‌کنش ترکیبات فعال پیشنهادی با گیرنده (شکل ۸-۳ تا شکل ۱۱-۳) نشان می‌دهد که این ترکیبات پیوندهای هیدروژنی مناسبی را با طول پیوند قابل قبول با اسید آمینه‌های کلیدی مانند Thr190، Glu166، His164، Phe140، Ser144، Leu141، Gln189 و Phe140 پیوندهای هیدروفوبی برقرار کرده‌اند. نتایج به‌دست آمده در خصوص برهم‌کنش‌های هیدروژنی و هیدروفوبی ترکیبات پیشنهادی با گیرنده، با نتایج گزارش شده در مقالات اخیر که داکینگ مولکولی بازدارنده‌های 3CL<sup>pro</sup> با گیرنده 6LU7 را گزارش داده‌اند، مطابقت دارد [۲۰۳-۲۰۵]. علاوه بر این، برخی از داروهای مرجع ضد ویروسی، مانند Lopinavir و Nelfinavir، که در برابر کووید-۱۹ مؤثر هستند، نیز در مطالعات داکینگ مولکولی مورد مطالعه قرار گرفته‌اند و برهم‌کنش‌های هیدروفوبی و هیدروژنی مشابهی را نشان داده‌اند [۲۰۵]. به‌طوری‌که پیوندهای هیدروژنی با اسید آمینه‌های کلیدی همچون Glu166، Arg188، His164، و Gln189 برقرار شده است و به‌علاوه پیوندهای هیدروفوبی با اسید آمینه‌های کلیدی همچون His163، Asp167، Ala191، Leu167، Met165، Leu141، Phe140، Ser114، و Gln189 گزارش شده است [۲۰۵، ۲۰۶]. علاوه بر این دو داروی مرجع، در اواسط اپیدمی کووید-۱۹، داروی Remdesivir نیز به‌عنوان یک مهارکننده مناسب در برابر SARS-CoV-2 اعلام عمومی شد. برهم‌کنش‌های هیدروژنی Remdesivir با اسید آمینه‌های کلیدی مانند Gly143، Thr24، His164، Glu166 و برهم‌کنش‌های هیدروفوبی نیز با اسید آمینه‌های کلیدی از جمله His41 و Met49 برقرار شده

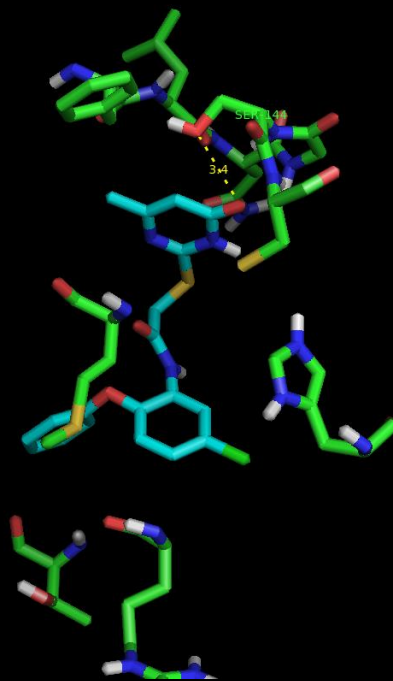
است [۲۰۷-۲۱۱]. با توجه به نتایج برهم کنش‌های مذکور در کمپلکس ترکیبات پیشنهادی - گیرنده نیز مشاهده می‌شود و فعال بودن ترکیبات پیشنهادی از حیث مطالعات داکینگ مولکولی نیز قابل توجه می‌باشد. علاوه بر مطالعه داکینگ مولکولی، پارامترهای فارماکوکینتیکی (PK) ترکیبات پیشنهادی و پارامتر سهولت سنتز این ترکیبات نیز با استفاده از ابزار وب رایگان Swiss-ADME محاسبه شد. پارامترهای PK محاسبه شده برای همه ترکیبات پیشنهادی در جدول ۳-۲ خلاصه شده است. نتایج نشان می‌دهد که همه ترکیبات پیشنهادی با پارامترهای قاعده پنج لیپینسکی (وزن مولکولی  $(MW < 500)$ )، چربی دوستی ( $4/15$ )  $(MLOGP < 10)$ ، تعداد گیرنده‌های پیوند هیدروژنی ( $\#H-B-acc < 10$ )، تعداد دهنده‌های پیوند هیدروژنی ( $\#H-B-don < 5$ ) و تعداد پیوندهای قابل چرخش ( $\#Rot-B < 10$ ) مطابقت دارند. علاوه بر این، فاکتور سهولت سنتز (Syn-Acc) هر ترکیب محاسبه و در جدول ۳-۲ آورده شده است. نتایج حاصله نشان می‌دهد که همه ترکیبات پیشنهادی دارای مقادیر فاکتور سهولت سنتز کم‌تر از ۱۰ هستند. در نتیجه، امکان سنتز آزمایشگاهی این ترکیبات وجود دارد و سنتز آن‌ها در آزمایشگاه با پیچیدگی خاصی همراه نیست. به‌طور خلاصه، نتایج حاصل از فعالیت مناسب پیش‌بینی شده مدل ALASSO-ANN، بر هم‌کنش‌های بالقوه به‌دست‌آمده از داکینگ مولکولی، پارامترهای قابل قبول قاعده پنج لیپینسکی و مقادیر سهولت سنتز مناسب، اهمیت سنتز ترکیبات پیشنهادی را به‌عنوان کاندیدهای مناسبی برای بازدارنده‌های جدید  $3CL^{PIV}$  را تأیید می‌نمایند.

DOI For evaluation only.  
tact sales@deisci.com.



۷۴

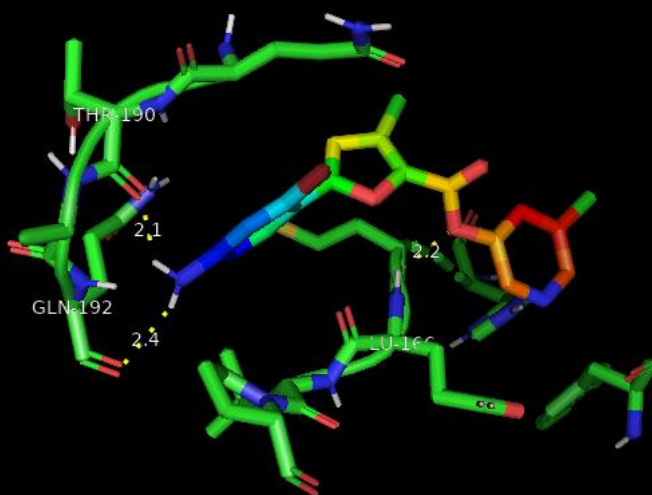
DOI For evaluation only.  
tact sales@deisci.com.



۷۲

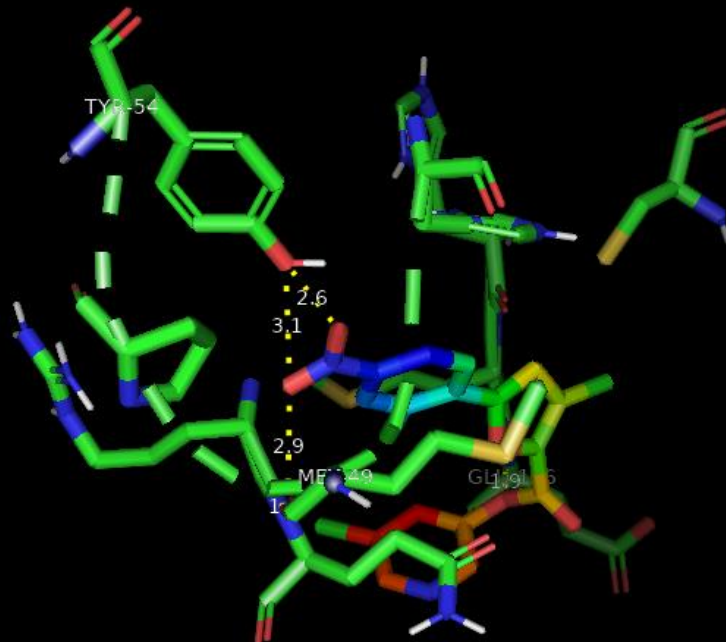
شکل ۷-۳ برهم کنش ترکیبات فعال (۷۴) و کم فعال (۷۲) موجود در مجموعه داده‌ها با اسید آمینه‌های کلیدی

No License File - For Evaluation Only (30 days remaining)



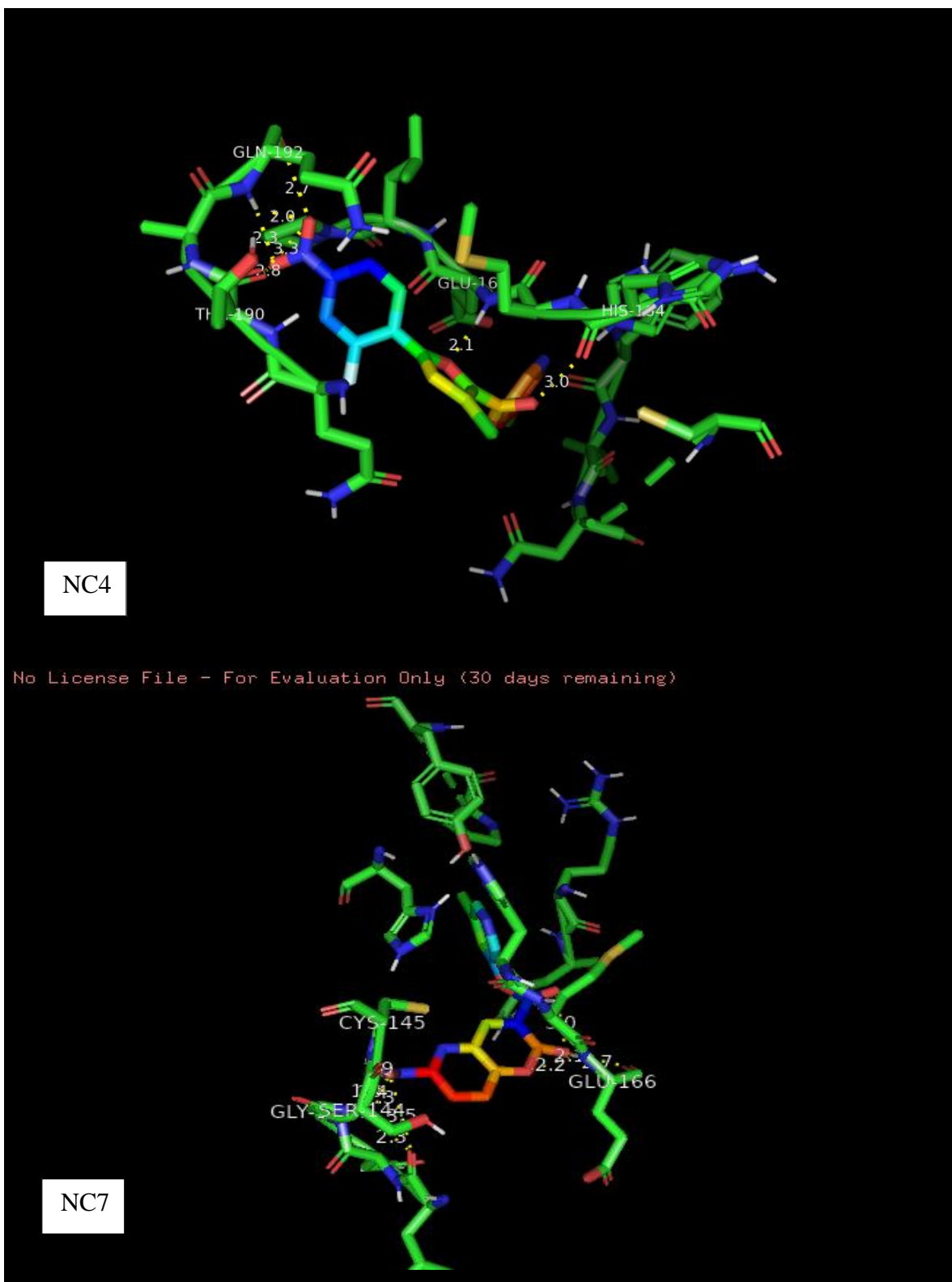
NC1

No License File - For Evaluation Only (30 days remaining)



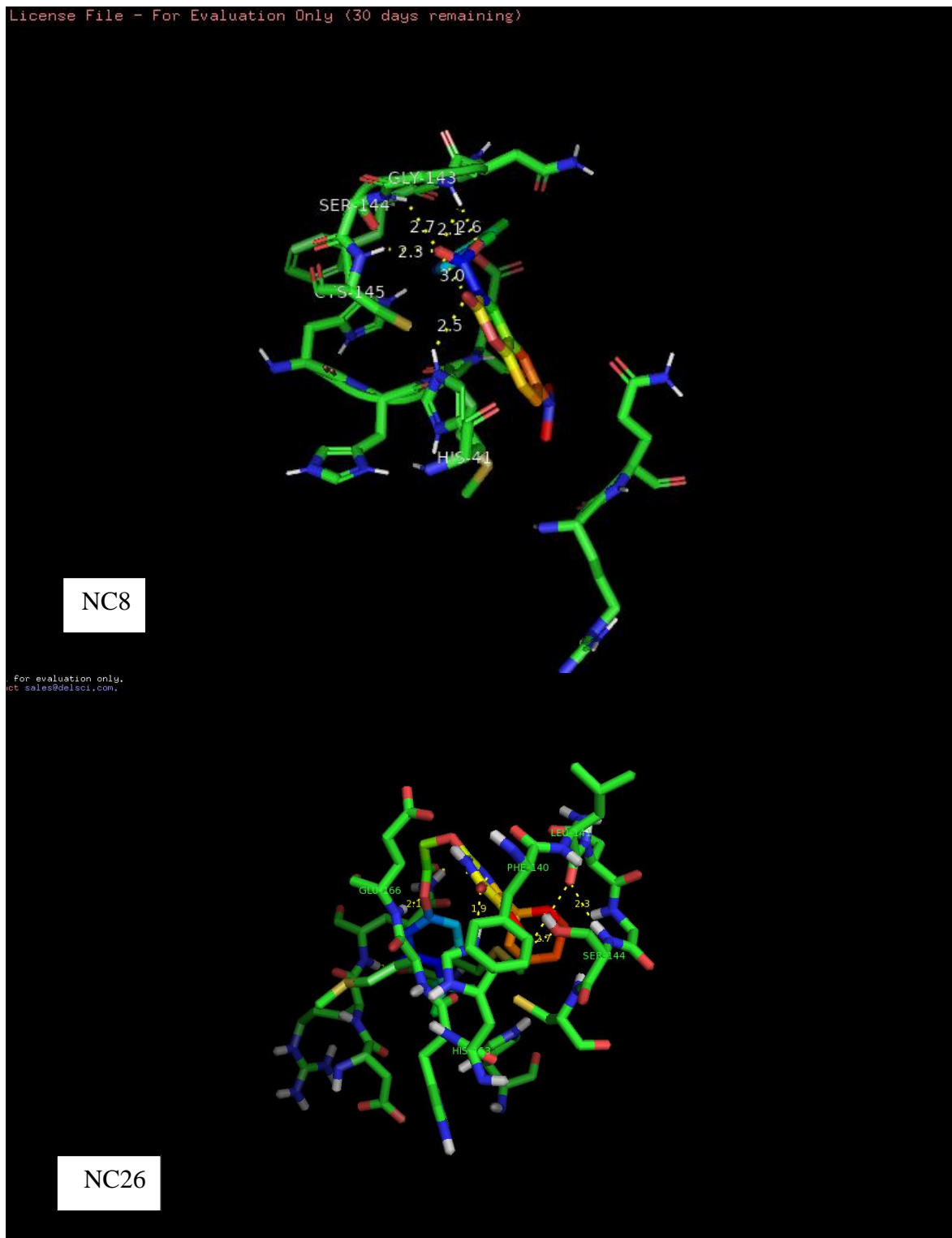
NC3

شکل ۸-۳ برهم کنش ترکیبات پیشنهادی (NC1 و NC3) با اسید آمینه‌های کلیدی



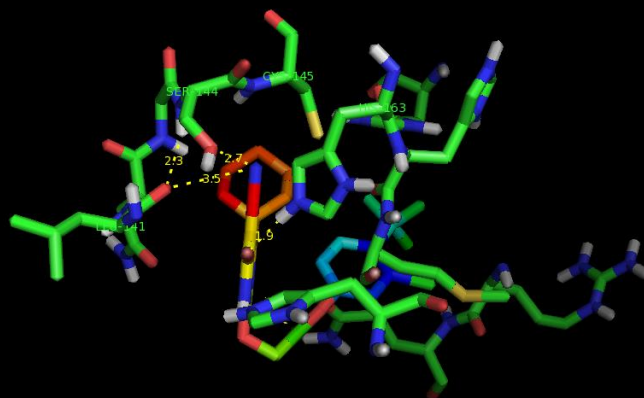
شکل ۹-۳ برهم کنش ترکیبات پیشنهادی (NC4 و NC7) با اسید آمینه‌های کلیدی





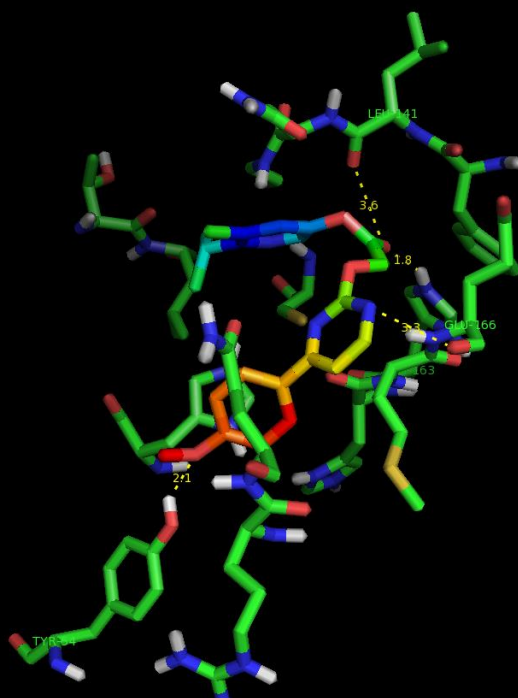
شکل ۱۰-۳ برهم کنش ترکیبات پیشنهادی (NC8 و NC26) با اسید آمینه‌های کلیدی

For evaluation only.  
get\_sales@deiscil.com.



NC28

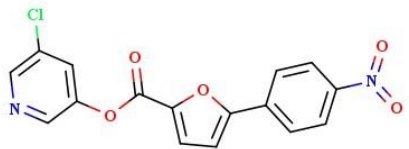
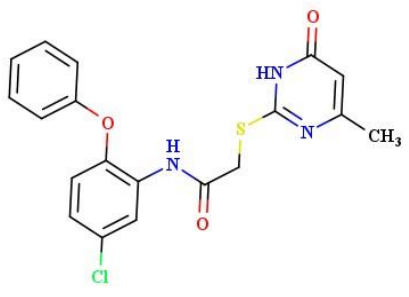
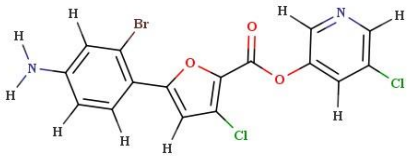
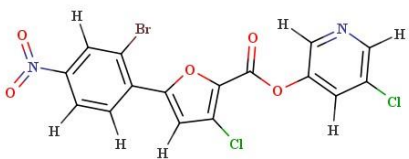
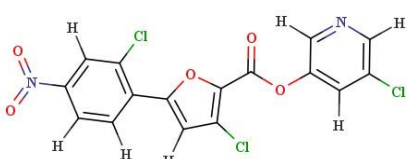
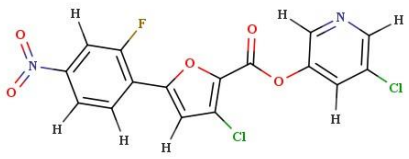
For evaluation only.  
get\_sales@deiscil.com.



NC30

شکل ۳-۱۱ برهم کنش ترکیبات پیشنهادی (NC28 و NC30) با اسید آمینه‌های کلیدی

جدول ۲-۳ پارامترهای PK محاسبه شده برای ترکیبات مورد مطالعه و ترکیبات پیشنهادی

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
۷۴		۳۴۴/۷۱	۱/۳۵	۵	۰	۶	۳/۰۳	۷/۲۲
۷۲		۴۰۱/۸۷	۲/۶۸	۷	۲	۴	۲/۸۵	۴/۰۰
NC1		۴۲۸/۰۶	۲/۹۴	۴	۱	۴	۳/۲	۶/۷۹
NC2		۴۵۸/۰۵	۲/۴۷	۵	۰	۶	۳/۲۳	۶/۷۹
NC3		۴۱۳/۶	۲/۳۶	۵	۰	۶	۳/۱۴	۶/۷۹
NC4		۴۱۵/۱۶	۱/۴۵	۵	۱	۸	۳/۲۱	۶/۷۹

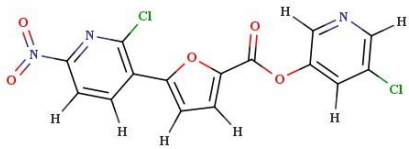
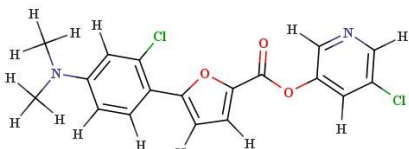
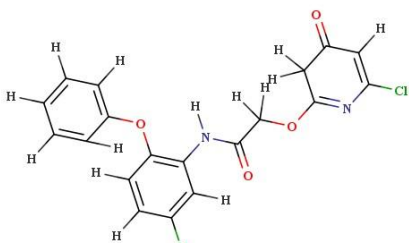
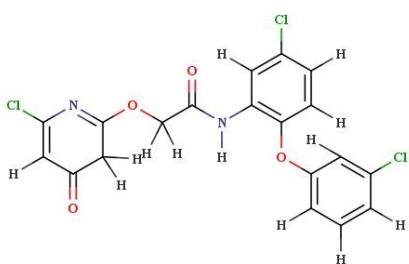
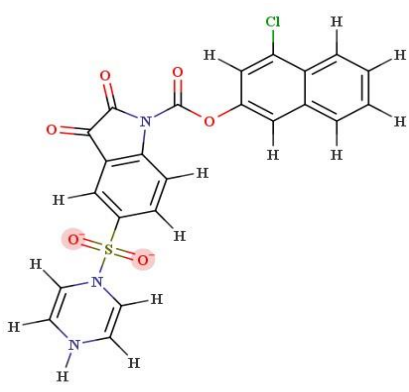
ادامه جدول ۲-۳

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
NC5		۳۱۷/۶۸	۱/۴	۳	۱	۶	۳/۰۱	۶/۷۹
NC6		۳۶۷/۱۶	۲/۷۱	۴	۱	۵	۳/۰۹	۶/۷۷
NC7		۳۴۷/۶۷	۱/۳۸	۴	۰	۸	۳/۰۲	۶/۷۲
NC8		۳۴۴/۷۵	۲/۲۵	۴	۰	۵	۳/۰۹	۶/۷
NC9		۳۴۹/۱۷	۲/۳۳	۴	۱	۴	۳/۰۹	۶/۶۸

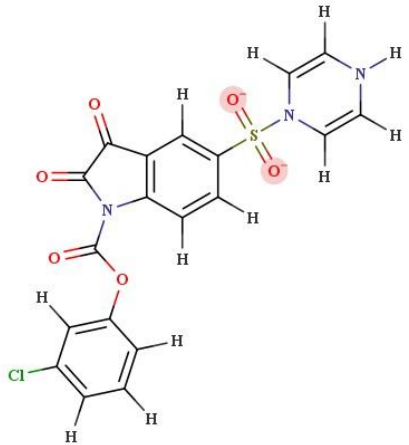
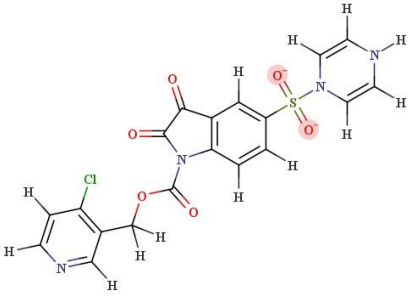
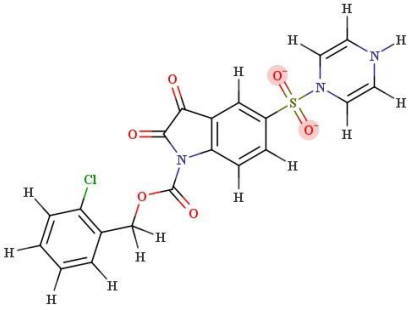
ادامه جدول ۲-۳

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
NC10		۳۱۴/۷۲	۱/۸۲	۴	۱	۴	۳/۰۱	۶/۶۶
NC11		۴۲۵/۶۱	۳/۱۲	۴	۰	۴	۳/۱۱	۶/۶۶
NC12		۳۷۸/۶	۳/۰۱	۴	۰	۴	۲/۹۹	۶/۶۵
NC13		۳۴۶/۶۸	۱/۳۳	۴	۰	۷	۲/۹۲	۶/۶۴
NC14		۳۷۸/۲۱	۲/۱۵	۵	۰	۵	۳/۳۶	۶/۶۴
NC15		۳۷۹/۱۵	۱/۸۶	۵	۰	۶	۳/۱	۶/۶۳

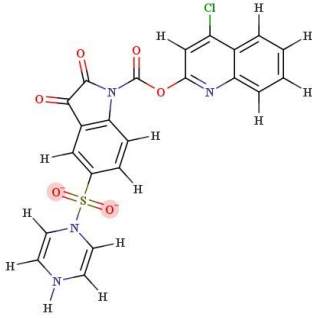
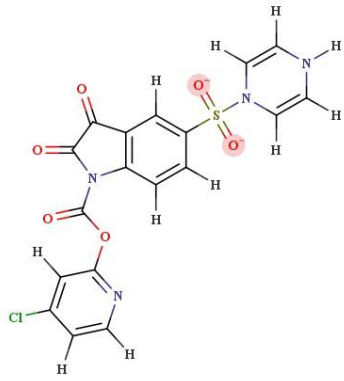
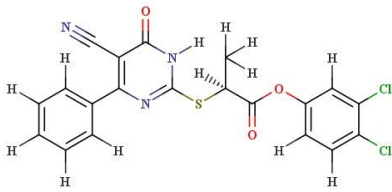
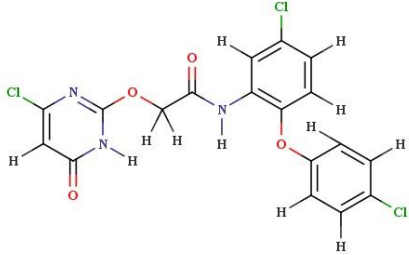
ادامه جدول ۲-۳

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
NC16		۳۸۰/۱۴	۱/۶۴	۵	۰	۷	۳/۲۱	۶/۵۷
NC17		۳۷۷/۲۲	۲/۷۹	۵	۰	۴	۳/۲۴	۶/۵۴
NC18		۴۰۵/۲۳	۲/۳	۷	۱	۵	۳/۵۹	۵/۹۲
NC19		۴۳۹/۶۸	۲/۷۹	۷	۱	۵	۳/۶	۵/۸۱
NC20		۴۹۵/۸۹	۲/۰۲	۵	۱	۶	۳/۳۳	۵/۷۶

ادامه جدول ۲-۳

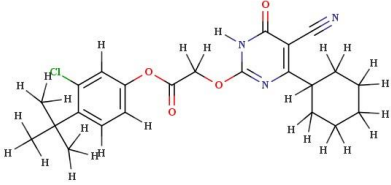
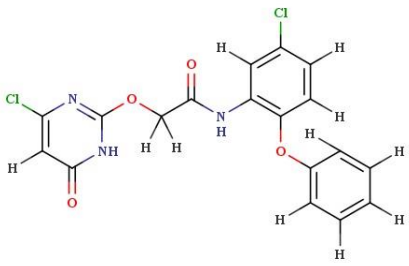
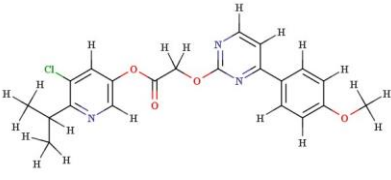
شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
NC21		۴۴۵/۸۳	۱/۷۳	۵	۱	۶	۳/۱۱	۵/۷۱
NC22		۴۶۰/۸۵	۰/۶۹	۶	۱	۷	۳/۱۷	۵/۶۴
NC23		۴۵۹/۸۶	۱/۶۸	۶	۱	۶	۳/۲۲	۵/۵۳

ادامه جدول ۲-۳

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
NC24		۴۹۶/۸۸	۱/۷۱	۵	۱	۷	۳/۳۶	۵/۴۵
NC25		۴۴۶/۸۲	۱/۱۴	۵	۱	۷	۳/۱۵	۵/۴۳
NC26		۴۴۶/۳۱	۲/۷۸	۶	۱	۵	۳/۴۹	۵/۴۰
NC27		۴۴۰/۶۶	۲/۶۴	۷	۲	۵	۲/۸	۵/۳۶



ادامه جدول ۲-۳

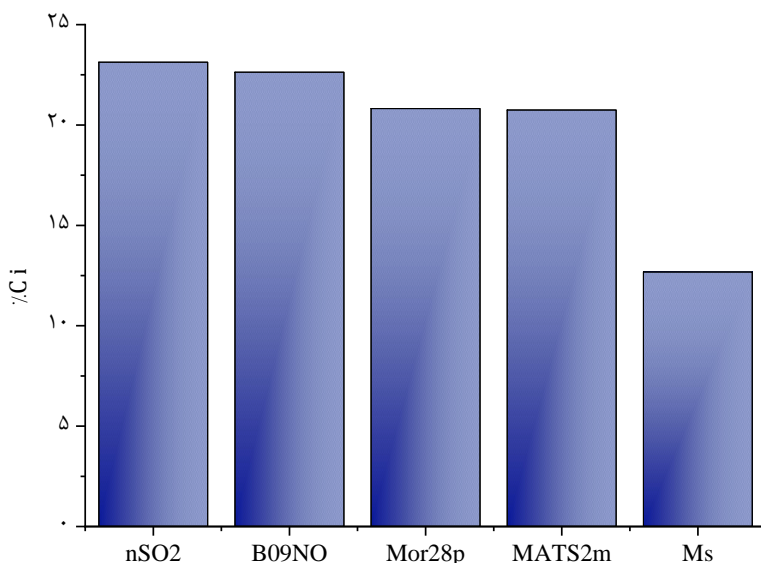
شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot- B	#H-B- don	#H-B- acc	Syn-Acc	pEC <sub>50</sub>
NC28		۴۴۳/۹۲	۲/۷۶	۷	۱	۶	۳/۶۸	۵/۳۳
NC29		۴۰۶/۲۲	۲/۱۵	۷	۲	۵	۲/۸	۵/۲۹
NC30		۴۱۳/۸۵	۲/۱۳	۸	۰	۷	۳/۲۷	۵/۲۶

### ۳-۱-۳ تجزیه و تحلیل توصیف‌کننده‌های مدل LAD-LASSO-LM-ANN برای

#### مجموعه داده‌های ضد ایدز و ضد سرطان

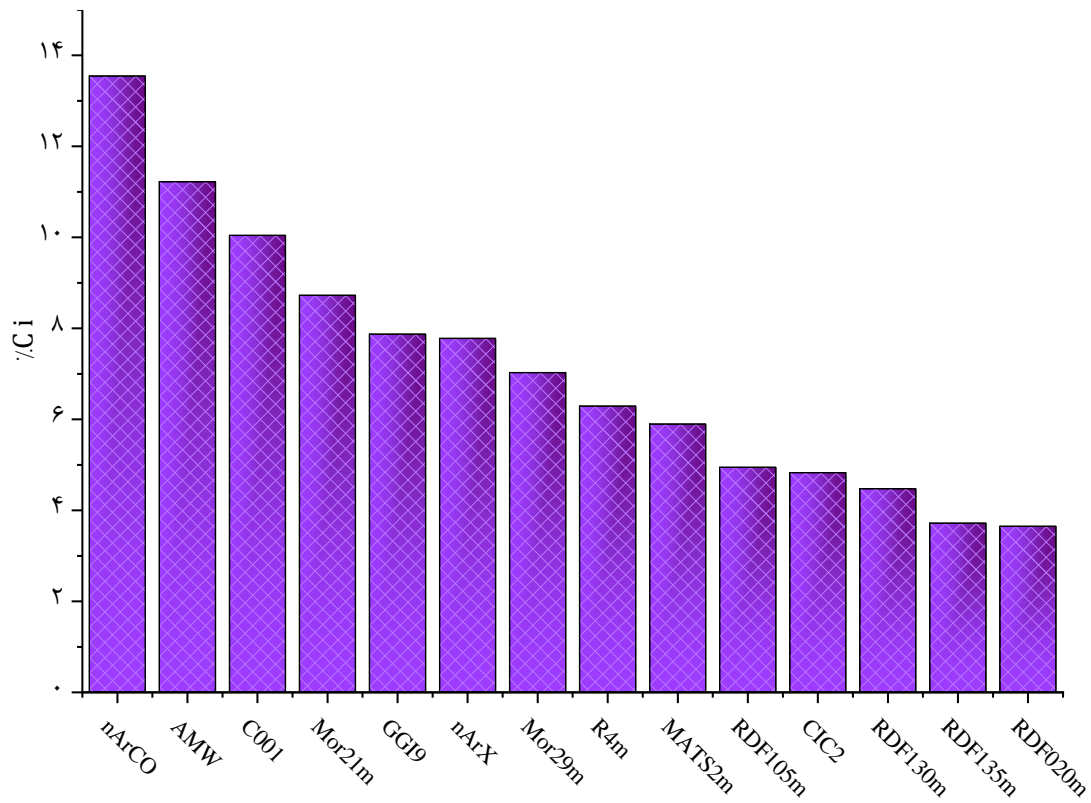
#### ۳-۱-۳-۱ محاسبه سهم مشارکت هر توصیف‌کننده در مدل LAD-LASSO-LM-ANN

مدل‌های QSAR برتر LAD-LASSO-LM-ANN برای مجموعه داده‌های مختلف با استفاده از توصیف‌کننده‌های منتخب (جدول ۲-۱۶) در فصل قبل ساخت و ارزیابی مدل بیان شد. میزان مشارکت هر توصیف‌کننده در مدل‌های LAD-LASSO-LM-ANN در شرایط بهینه برای مجموعه داده‌های مختلف، با محاسبه درصد سهم مشارکت هر توصیف‌کننده ( $\%C_i$ ) مطابق با توضیحات مندرج در بخش ۱-۵-۸-۷ و بر اساس رابطه ۱-۱۵ برآورد شد. نتایج به‌دست آمده برای سهم مشارکت توصیف‌کننده‌های منتخب هر کدام از مجموعه داده‌ها در مدل LAD-LASSO-LM-ANN مربوطه در شکل ۳-۱۲ تا شکل ۳-۱۴ آورده شده است. در ادامه، با توجه به درصد سهم مشارکت توصیف‌کننده‌های منتخب در مدل‌های بهینه، به توضیح برخی از مهم‌ترین توصیف‌کننده‌ها پرداخته خواهد شد و کاربرد مدل‌های ارائه شده در پیشنهاد ترکیبات جدید با فعالیت دارویی مناسب بیان می‌گردد.

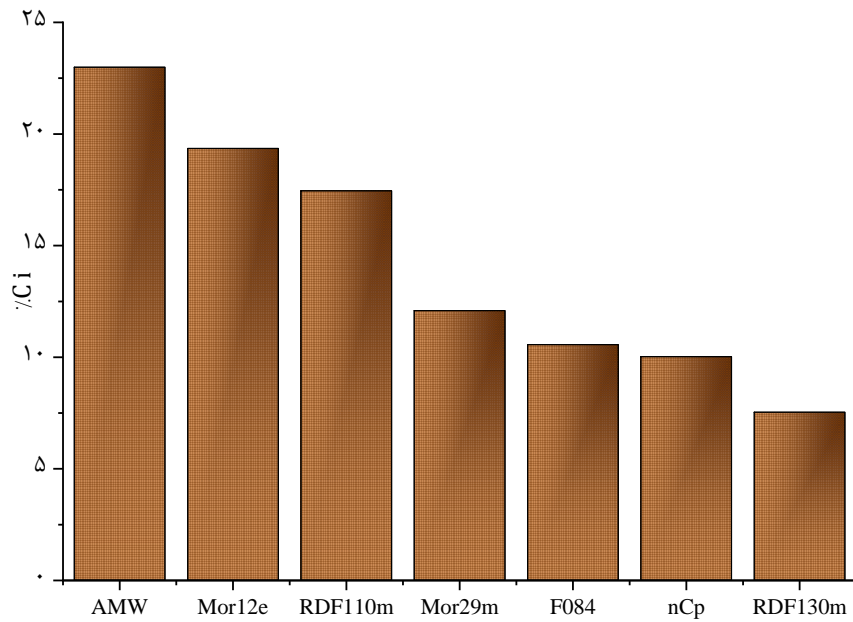


شکل ۳-۱۲ نام توصیف‌کننده‌های منتخب روش LAD-LASSO برای مجموعه داده‌های ایدز

شکل ۳-۱۲ نمودار سهم مشارکت توصیف‌کننده‌ها در مدل LAD-LASSO-LM-ANN برای بازدارنده‌های ضد ایدز



شکل ۳-۱۳ نمودار سهم مشارکت توصیف کننده‌ها در مدل LAD-LASSO-LM-ANN برای بازدارنده‌های سرطان کارسینوم کولورکتال نام تو و صیف کننده‌های منتخب روش LAD-LASSO به رای مجموعه بازدارنده‌های سرطان کولورکتال



شکل ۳-۱۴ نمودار سهم مشارکت توصیف کننده‌ها در مدل LAD-LASSO-LM-ANN برای بازدارنده‌های سرطان ریه نام تو و صیف کننده‌های منتخب روش LAD-LASSO به رای مجموعه بازدارنده‌های ریه

۳-۱-۲ بررسی میزان و چگونگی تأثیر توصیف کننده‌ها بر فعالیت دارویی ترکیبات مجموعه داده‌های ضد ایدز و کاربرد مدل LAD-LASSO-LM-ANN ارائه شده در پیشنهاد

### ترکیبات جدید

برای بررسی تأثیر هر توصیف کننده بر معادله نهایی، مدل LAD-LASSO-LM-ANN برای مجموعه داده‌های ضد ایدز استخراج شد. تأثیر افزایش و کاهش مقادیر هر توصیف کننده بر فعالیت دارویی ترکیبات با توجه به ضرایب استاندارد شده توصیف کننده‌ها (جدول ۲-۱۶) مورد بررسی قرار گرفت. در ادامه شرح مختصری از برخی از توصیف کننده‌ها، با توجه به نتایج سهم مشارکت و تأثیر ضرایب توصیف کننده‌های منتخب مدل‌های LAD-LASSO، آورده شده است. در نهایت برهم کنش ترکیبات پیشنهادی با گیرنده با استفاده از مطالعه داکینگ مولکولی مورد ارزیابی قرار گرفت.

از بین توصیف کننده‌های مدل LAD-LASSO-LM-ANN برای بازدارنده‌های ایدز، توصیف کننده  $nSO_2$  مربوط به وجود گروه‌های  $SO_2$  در ترکیبات مورد مطالعه است. ضریب اثر مثبت این توصیف کننده در مدل LAD-LASSO-LM-ANN (جدول ۲-۱۶) نشان‌دهنده این است که حضور این گروه در ترکیب موجب افزایش فعالیت دارویی ترکیب می‌شود. برای مثال دو ترکیب شماره ۱۴ و ۴۰ از نظر اثر وجود و عدم وجود گروه  $SO_2$  مورد مقایسه قرار گرفتند. ترکیب (ساختار شماره ۱۴) با داشتن گروه  $SO_2$  دارای فعالیت دارویی بیش‌تری ( $pIC_{50} = ۸/۲۶$ ) نسبت به ترکیب فاقد گروه  $SO_2$  (ساختار شماره ۴۰ با  $pIC_{50} = ۴/۵۲$ ) دارد. توصیف کننده بعدی  $B09[N-O]$  است. این توصیف کننده مربوط به وجود/عدم وجود  $N-O$  در فاصله مکانی  $۹\text{\AA}$  از مرکز است. ضریب اثر منفی این توصیف کننده در مدل LAD-LASSO-LM-ANN نشان‌دهنده این است که در نبود  $N-O$  در فاصله مکانی  $۹\text{\AA}$ ، ترکیبات مورد مطالعه فعالیت دارویی بیش‌تری دارند. این توصیف کننده بیان ارتباط ساختار-فعالیت ترکیباتی با گروه‌های نیترو و سیانو در مجاورت نیتروپیریدین، بنزونیتریل، پیریدازین و پیریدین در مدل ظاهر شده است. به‌عنوان مثال ترکیبات شماره

۳۶، ۵۶، ۶۴ و ۶۹ به دلیل مجاورت  $9 \text{ \AA}$  با گروه نیترو و سیانو با حلقه‌های مذکور دارای فعالیت دارویی به ترتیب برابر با  $4/79$ ،  $5/29$ ،  $5/55$  و  $5/85$  هستند و از دسته ترکیباتی با فعالیت دارویی ضعیف و متوسط به شمار می‌آیند. توصیف کننده Mor28p مربوط به طبقه 3D-MoRSE است [۹، ۲۱۲]. این دسته از این توصیف کننده‌ها به گونه‌ای تعریف می‌شوند که سهم ویژگی اتمی قطبش پذیری (p)، را در یک زاویه پراکندگی تعیین شده به ویژگی هدف منعکس می‌کنند و امکان تمایز بین ماهیت اتم را فراهم می‌کند. Mor28p، با زاویه پراکندگی  $28 \text{ \AA}^{-1}$ ، توسط قطبش پذیری اتمی وزن دهی می‌شود. این توصیف کننده دارای مقادیر منفی می‌باشد و با توجه به ضریب منفی Mor28p، نشان می‌دهد که وجود برخی از توده‌های اتمی مطلوب سبب بهبود فعالیت بازدارنده‌های ایدز می‌شود [۲۱۳]. برای مثال دو ترکیب شماره ۲۰ و ۳۶ مورد مقایسه قرار گرفت. همان‌طور که گفته شد، مطابق با ضریب اثر منفی توصیف کننده بر فعالیت دارویی، ترکیب (ساختاره شماره ۲۰) با منفی‌ترین مقدار Mor28p دارای فعالیت دارویی بیش‌تری نسبت به ترکیب (ساختار شماره ۳۶) با مثبت‌ترین مقدار Mor28p است. توصیف کننده مهم دیگر MATS2m است. این توصیف کننده مربوط به طبقه‌بندی خودهمبستگی دوبعدی وزن‌دهی شده با جرم مولکولی است و اطلاعاتی در مورد توزیع جرم مولکولی در امتداد ساختار توپولوژیکی ارائه می‌دهد [۲۱۴]. این توصیف کننده با ضریب اثر مثبت، باعث افزایش فعالیت دارویی می‌شود.

Ms توصیف کننده نشان‌دهنده میانگین حالت الکتروتوپولوژیکی و از دسته توصیف کننده‌های شاخص‌های اساسی است. Ms دارای علامت مثبت در مدل LAD-LASSO است. بنابراین ترکیباتی که دارای مقدار بزرگ‌تری از این توصیف کننده هستند فعالیت دارویی بیش‌تری را به خود اختصاص می‌دهند. در مجموعه داده‌های بازدارنده‌های ایدز، ترکیباتی که بر روی حلقه فنیل بال راست دارای گروه‌های فلئور در مجموع فعالیت دارویی بیشتری نسبت به ترکیبات مشابه ولی فاقد گروه فلئور هستند. برای مثال ترکیب شماره ۳۹ (فعالیت دارویی برابر با  $6/14$ ) دارای سه گروه فلئور بر روی حلقه فنیل می‌باشد و در

مقابل ترکیب شماره ۳۳ (فعالیت دارویی برابر با ۵/۳۷) درای دو گروه متوکسی بر روی حلقه فنیل می باشد. از این رو بازدارندگی ترکیب ۳۵ کم تر از ترکیب ۳۹ می باشد.

یکی از مهم ترین کاربردهای مدل QSAR، پیشنهاد ترکیبات جدید با استفاده از توصیف کننده های مدل برتر پیشنهادی است. از این رو با توجه به اثر هر توصیف کننده بر فعالیت دارویی مربوطه، می توان به چگونگی ایجاد یک اصلاح ساختاری مؤثر که منجر به افزایش فعالیت دارویی شود، پی برد.

بنابراین با توجه به توضیح های ارائه شده در بخش ۱-۳-۱-۳ برای مدل QSAR بازدارنده های ایدز، تلاش شد تا در راستای تأثیر توصیف کننده ها بر فعالیت دارویی و انجام برخی از اصلاحات ساختاری بر روی ترکیبات ضعیف مورد مطالعه (ترکیبات ۴۰ و ۳۶ با  $pIC_{50}$  به ترتیب برابر با ۴/۵۲ و ۴/۷۷)، ترکیباتی با فعالیت دارویی مناسب پیشنهاد شود. از این رو با توجه به تأثیر وجود گروه  $SO_2$  و نبود گروه  $N-O$  در فاصله توپولوژیکی  $9\text{\AA}$ ، ترکیبات جدیدی با فعالیت دارویی مناسب در محدوده ترکیبات فعال، پیشنهاد شد و جزییات ساختاری آن ها در جدول ۳-۳ آورده شده است.

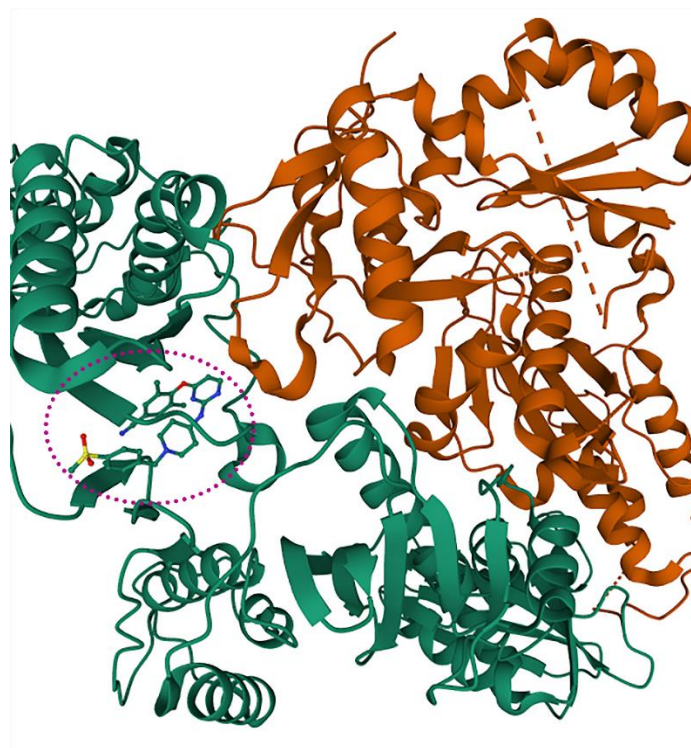
علاوه بر این ترکیبات جدیدی، با توجه به اثر توصیف کننده های مدل بهینه LAD-LASSO-LM-ANN بر فعالیت دارویی بازدارنده های سرطان کارسینوم کولورکتال، پیشنهاد شد. از این رو ترکیباتی با تعداد کتون های آروماتیک کم تر، تعداد گروه های هالوژنه متصل به حلقه آروماتیک بیش تر و تعداد اتم هایی با وزن اتمی بالاتر ایجاد شدند. جزییات ساختاری این ترکیبات پیشنهادی با فعالیت دارویی پیش بینی شده در جدول ۴-۳ آورده شده است. در نهایت با توجه به شرح مختصر توصیف کننده های منتخب مدل LAD-LASSO-LM-ANN برای بازدارنده های سرطان ریه، استنباط می شود که وزن مولکولی متوسط (وزن مولکولی/تعداد کل اتم ها) بالای ترکیبات، حضور گروه الکترون گاتیو از جمله هالوژن ها در ترکیبات، وجود گروه کربن  $Sp^3$  انتهایی و عدم وجود گروه  $F$  متصل به کربن  $Sp^2$  در ساختار ترکیب باعث افزایش فعالیت دارویی می شود. بنابراین با انجام برخی اصلاحات، ترکیبات جدیدی با فعالیت دارویی مناسب پیشنهاد شد

که جزییات ساختاری این ترکیبات پیشنهادی با فعالیت دارویی قابل قبول در جدول ۳-۴ آورده شده است. پس از اینکه ترکیبات مورد نظر پیشنهاد شدند ساختار هریک با استفاده از برنامه هایپرکم رسم و بهینه سازی شد. توصیف کننده های منتخب هر کدام از مدل های ANN -LAD-LASSO-LM-LM برای مجموعه داده های مربوطه با استفاده از نرم افزار دراگون محاسبه شدند. مقادیر فعالیت دارویی ترکیبات جدید با استفاده از مدل در شرایط بهینه پیش بینی شد و نتایج مربوط به فعالیت های دارویی پیش بینی شده ترکیبات جدید هر کدام از مجموعه داده ها، در جدول ۳-۴ آورده شد. در ادامه بحث، برای اثبات فعالیت ترکیبات پیشنهاد شده، از مطالعه داکینگ مولکولی استفاده شد که در بخش ۳-۱-۳ به طور کامل به آن پرداخته خواهد شد.

### ۳-۱-۳ مطالعه داکینگ مولکولی بازدارنده های ایدز

پس از ایجاد ترکیبات فعال پیشنهادی با فعالیت دارویی قابل قبول، برهم کنش های ترکیبات پیشنهادی با اسید آمینه های کلیدی جایگاه فعال گیرنده مورد نظر نیز مورد بررسی قرار گرفت. به طوری که ابتدا ساختار کریستالوگرافی (3M8Q) به پیشنهاد مقالات منتشر شده از سایت بانک اطلاعاتی پروتئین با فرمت pdb دانلود شد [۱۷۰, ۱۷۲]. ساختار کریستالوگرافی 3M8Q با ارزش تفکیک برابر با  $2/70^\circ A$  برای انجام محاسبات داکینگ مولکولی از کیفیت مناسبی برخوردار است [۲۱۵]. شکل ۳-۱۵ زنجیره اسید آمینه ای ساختار کریستالوگرافی 3M8Q و لیگاند کریستالوگرافی موجود در زنجیره را نشان می دهد. جایگاه فعال مربوط به ساختار کریستالوگرافی بر اساس مختصات مرکز ثقل لیگاند کریستالوگرافی تعریف شد. فرایند اعتبار سنجی داکینگ مولکولی با استفاده از زنجیره اسید آمینه ای A و لیگاند کریستالوگرافی انجام شد. نتایج داکینگ، تعداد اجرای الگوریتم ژنتیک برابر با ۱۵۰ را به عنوان داکینگ بهینه مشخص کرد. بنابراین داکینگ ترکیبات جدید پیشنهادی با ساختار کریستالوگرافی مربوط به هر مجموعه داده در تعداد اجرای الگوریتم ژنتیک برابر با ۱۵۰ انجام شد. از این رو، ساختارهای بهینه شده ترکیبات مورد مطالعه و

پیشنهادی با پسوند pdb ذخیره شدند و به‌عنوان ورودی لیگاند نرم‌افزار داکینگ مولکولی تعریف شدند. سپس داکینگ ترکیبات فعال و غیر فعال پیشنهادی با ساختار کریستالوگرافی مربوطه در شرایط بهینه داکینگ انجام شد. در نهایت بهترین پی‌کربندی مربوط به ترکیبات با کم‌ترین انرژی اتصال از نرم‌افزار داکینگ استخراج شدند و برهم‌کنش‌های لیگاند-گیرنده با استفاده از نرم‌افزار Discovery Studio Visualizer به‌دست آمد.

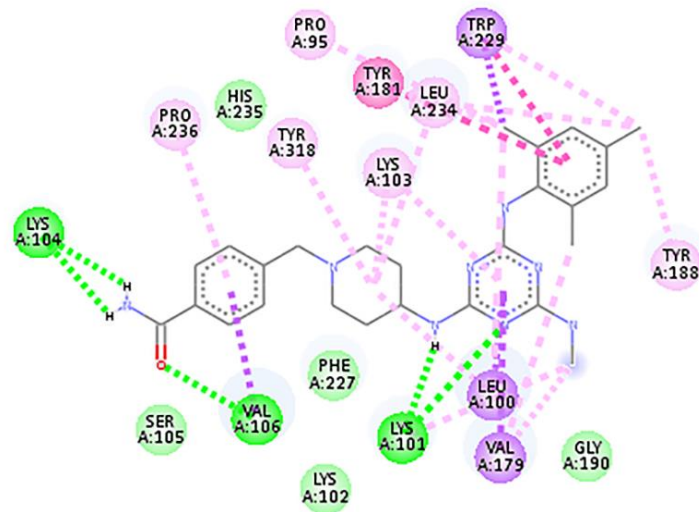


شکل ۳-۱۵ ساختار کریستالوگرافی 3M8Q [۲۱۵] (منطقه نقطه چین نشان‌دهنده لیگاند کریستالوگرافی و مابقی زنجیره‌های اسید آمینه‌ای است)

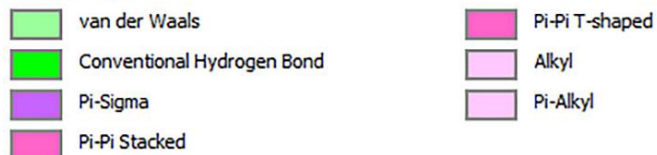
برهم‌کنش‌های متفاوت ترکیبات فعال و کم‌فعال با گیرنده در شکل ۳-۱۶ تا شکل ۳-۱۸ آورده شده است. با توجه به نتایج مشاهده می‌شود که ترکیب فعال شماره ۲۰ با  $pIC_{50}$  برابر با ۸/۳۴ با اسید آمینه‌های متفاوتی از جمله Lys104، Val106 و Lys101 برهم‌کنش‌های هیدروژنی برقرار کرده است، بنابراین با توجه به توانایی ترکیب در برقراری پیوندهای هیدروژنی و هیدروفوبی متفاوت، فعالیت این ترکیب نیز اثبات می‌شود. علاوه بر این به‌منظور مقایسه این ترکیب با ترکیبات ضعیف مورد مطالعه، برهم‌کنش‌های



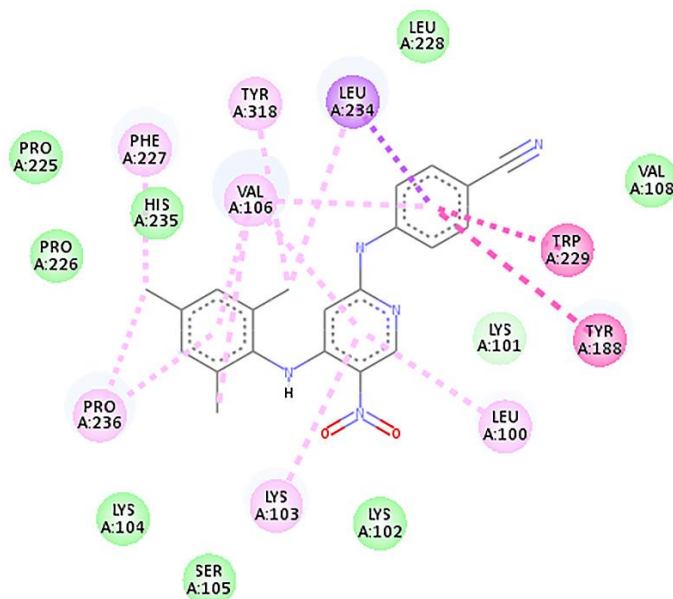
ترکیبات ضعیف ۴۰ و ۳۶ نیز مورد بررسی قرار گرفت و به ترتیب در شکل ۳-۱۷ و شکل ۳-۱۸ آورده شد. نتایج نشان می‌دهد که این دو ترکیب در برقراری پیوند هیدروژنی موفق نبوده‌اند. الگوی مربوط به برهم کنش ترکیبات مورد مطالعه برای مقایسه فعالیت ترکیبات جدید مد نظر قرار گرفتند. برهم کنش برخی از ترکیبات از جمله NC1، NC2، NC4 و NC6 در شکل ۳-۱۹ تا شکل ۳-۲۴ نشان داده شده است. با توجه به نتایج مشاهده می‌شود که این ترکیبات دارای غلظت مؤثر کم‌تر از ۱ میکرومولار هستند و در محدوده ترکیبات فعال قرار دارند. علاوه بر این که فعالیت این ترکیبات با استفاده از پیش‌بینی مدل LAD-LASSO- LM-ANN مربوطه تأیید شد، در ادامه برهم کنش‌های این ترکیبات نیز مورد بررسی قرار گرفت. به طوری که ترکیب NC1 به عنوان یکی از فعال‌ترین ترکیبات پیشنهادی، برهم کنش‌های هیدروژنی مناسبی با اسید آمینه‌های کلیدی Lys101، Val106، His235، Pro236، Tyr188 و Val179 برقرار کرده است که این نتیجه گواه مناسبی بر فعالیت چشم‌گیر این ترکیب می‌باشد. ترکیبات NC2، NC4 و NC6 نیز به ترتیب پیوندهای هیدروژنی مناسبی را با اسید آمینه‌های Leu234، Pro236، Lys103، Lys101 و Val106 برقرار کرده‌اند.



**Interactions**



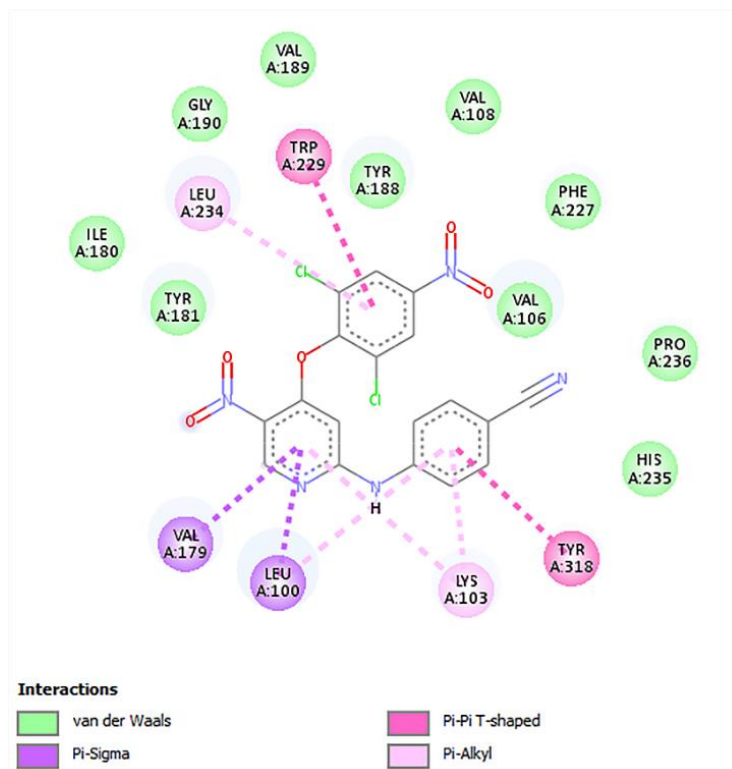
شکل ۱۶-۳ برهم کنش ترکیب فعال (۲۰) موجود در مجموعه داده‌های ضد ایدز با اسید آمینه‌های کلیدی



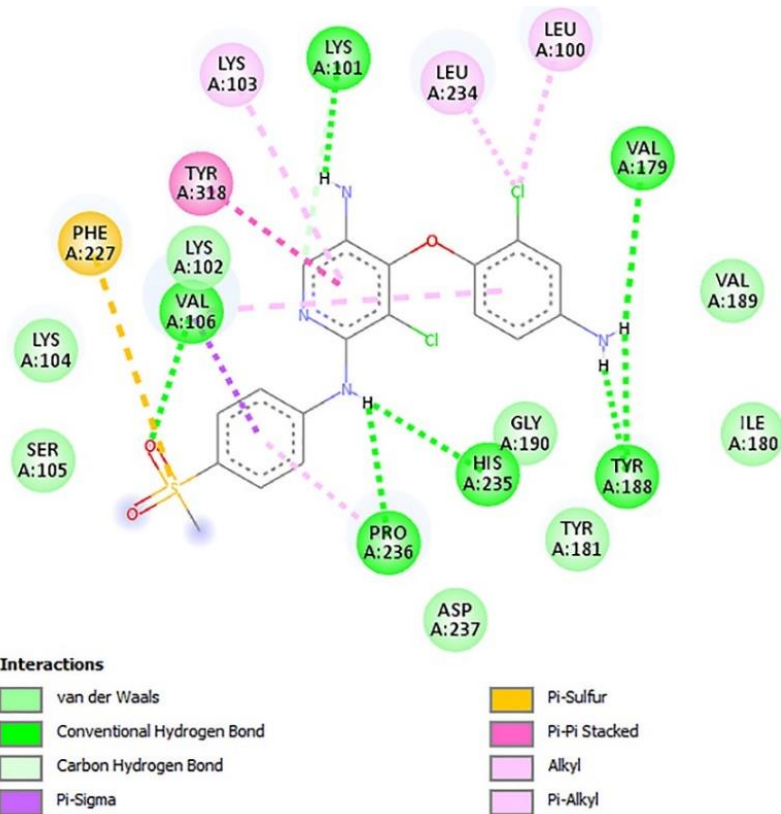
**Interactions**



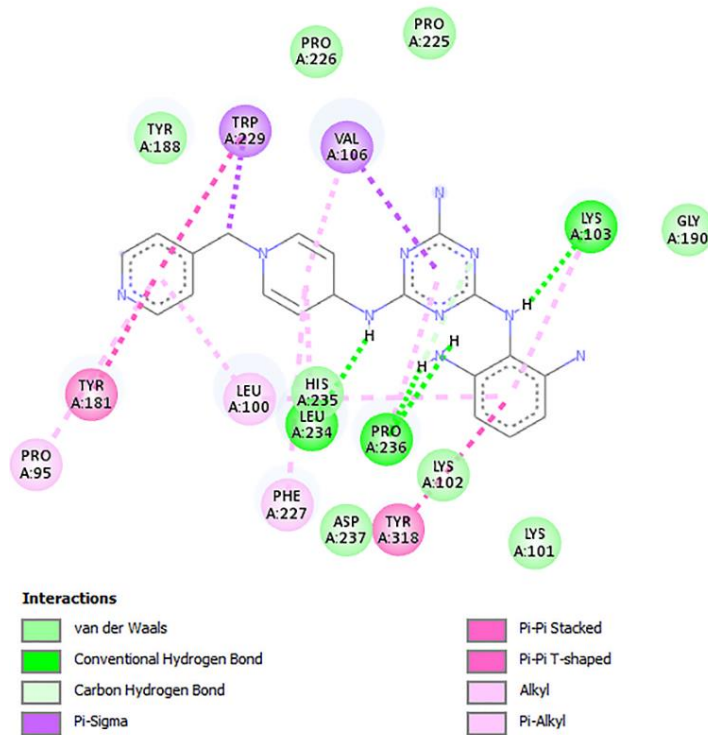
شکل ۱۷-۳ برهم کنش ترکیب کم فعال (۴۰) موجود در مجموعه داده‌های ضد ایدز با اسید آمینه‌های کلیدی



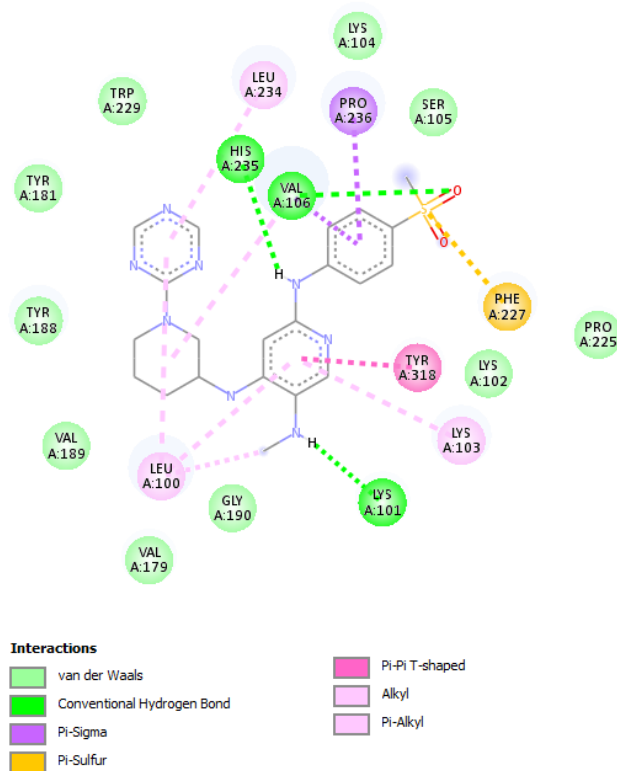
شکل ۳-۱۸ برهم کنش ترکیب کم فعال (۳۶) موجود در مجموعه داده‌های ضد ایدز با اسید آمینه‌های کلیدی



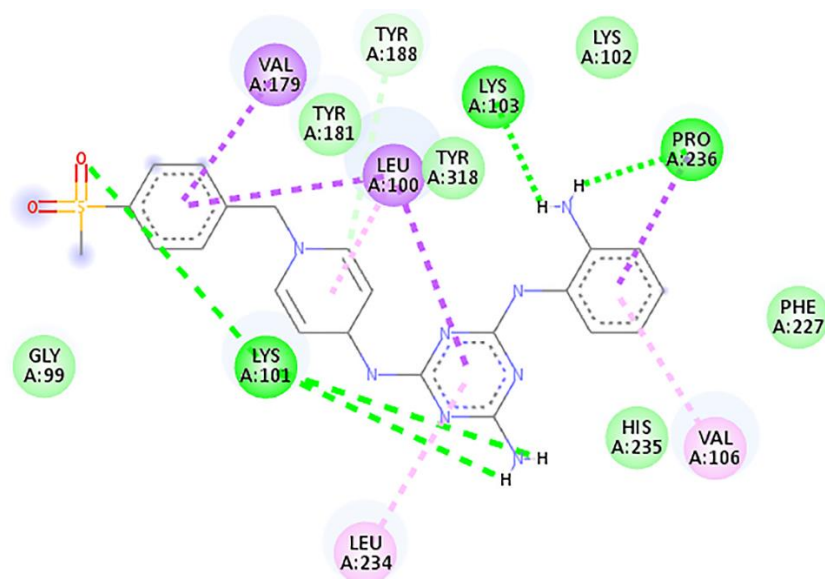
شکل ۳-۱۹ برهم کنش ترکیب پیشنهادی NCI با اسید آمینه‌های کلیدی



شکل ۲۰-۳ برهم کنش ترکیب پیشنهادی NC2 با اسید آمینه‌های کلیدی



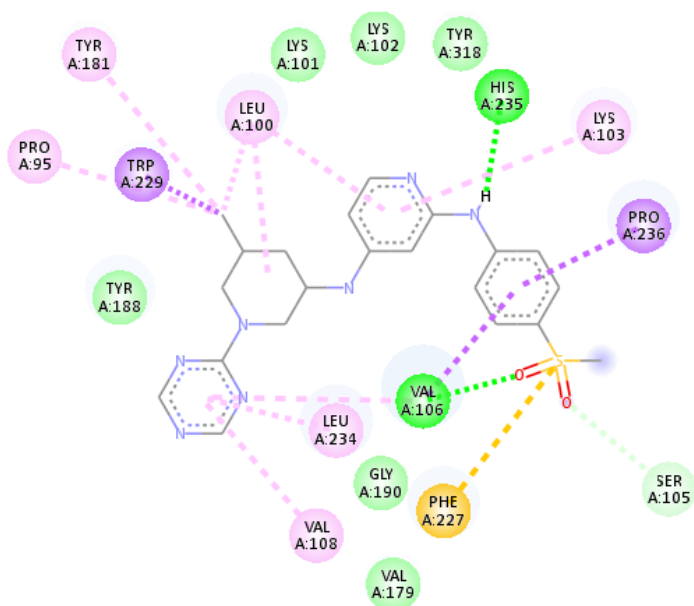
شکل ۲۱-۳ برهم کنش ترکیب پیشنهادی NC3 با اسید آمینه‌های کلیدی



**Interactions**

van der Waals	Pi-Sigma
Conventional Hydrogen Bond	Alkyl
Carbon Hydrogen Bond	Pi-Alkyl

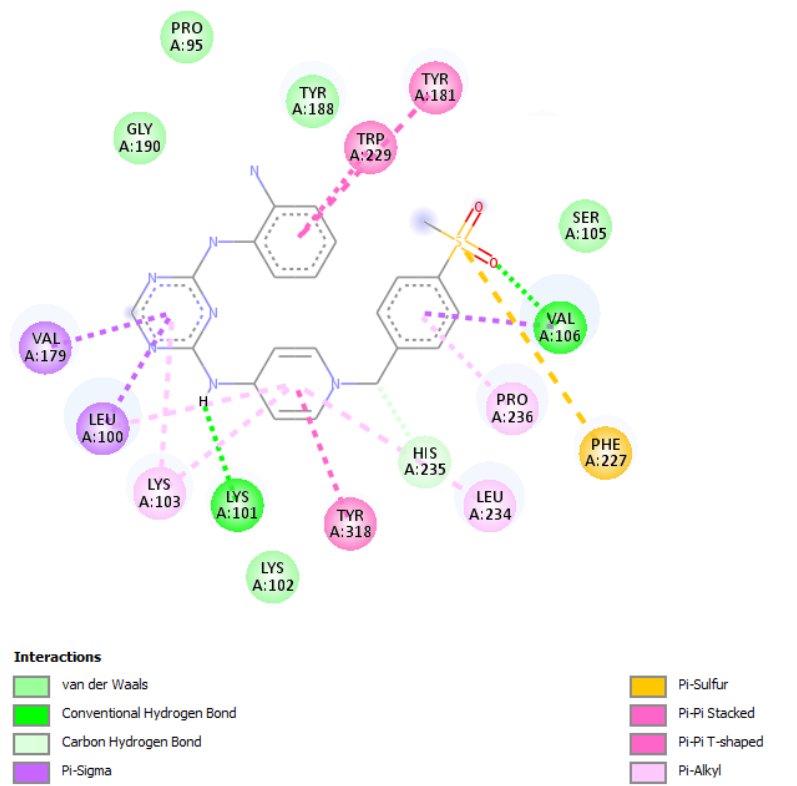
شکل ۲۲-۳ برهم کنش ترکیب پیشنهادی NC4 با اسید آمینه‌های کلیدی



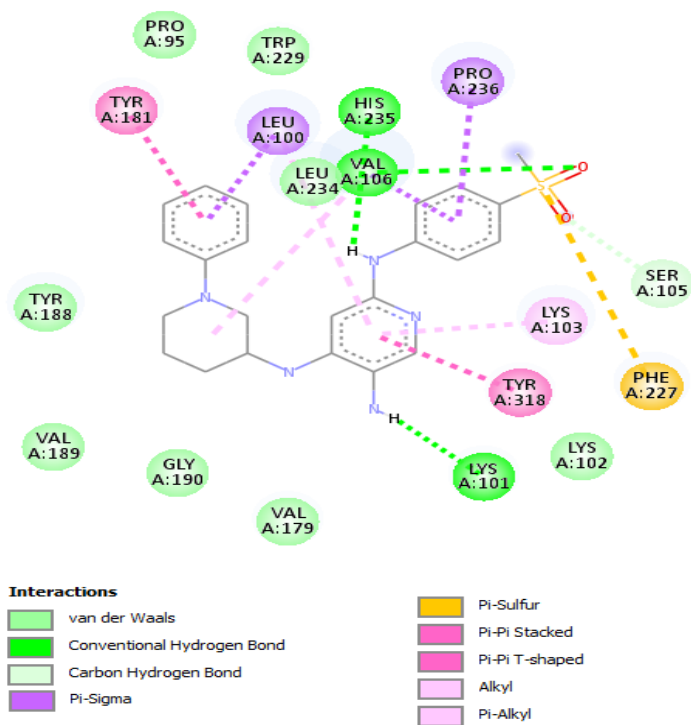
**Interactions**

van der Waals	Pi-Sulfur
Conventional Hydrogen Bond	Alkyl
Carbon Hydrogen Bond	Pi-Alkyl
Pi-Sigma	

شکل ۲۳-۳ برهم کنش ترکیب پیشنهادی NC5 با اسید آمینه‌های کلیدی

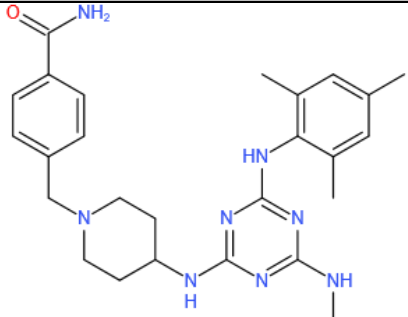
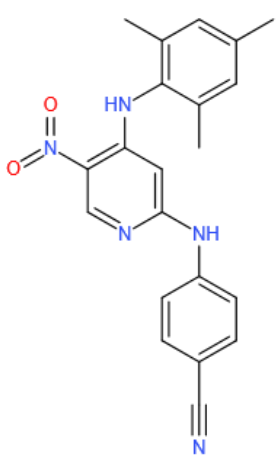
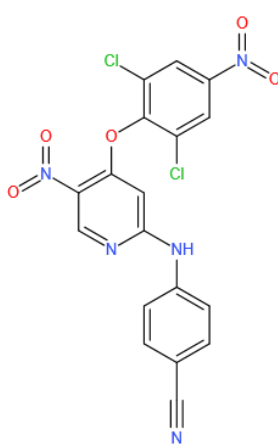


شکل ۲۴-۳ برهم کنش ترکیب پیشنهادی NC6 با اسید آمینه‌های کلیدی

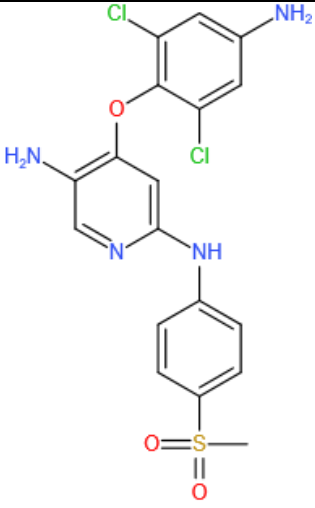
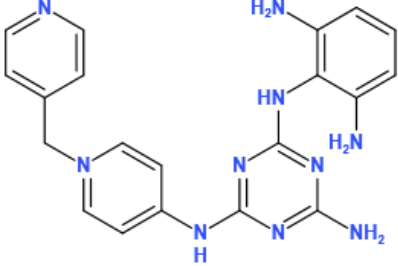
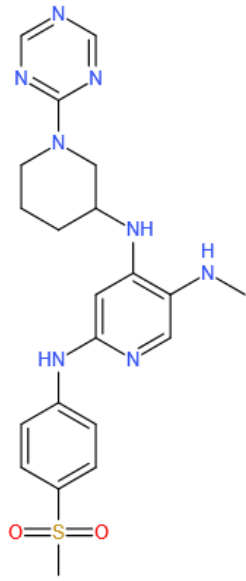


شکل ۲۵-۳ برهم کنش ترکیب پیشنهادی NC7 با اسید آمینه‌های کلیدی

جدول ۳-۳ پارامترهای PK محاسبه شده برای ترکیبات مورد مطالعه ضد ایدز و ترکیبات پیشنهادی

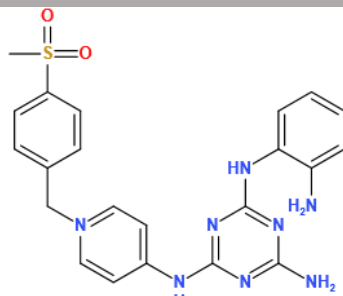
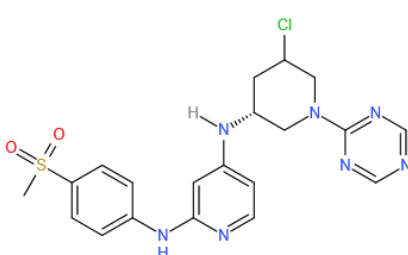
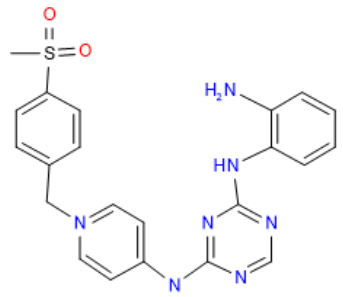
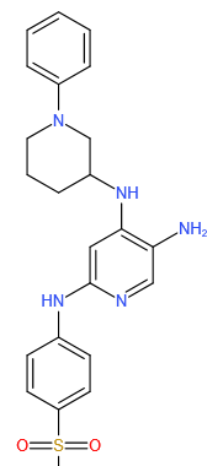
شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
۲۰		۴۷۴/۶	۱/۶۶	۸	۴	۵	۳/۷۸	۸/۳۴
۴۰		۳۷۳/۴۱	۲/۱	۵	۲	۴	۳/۳۵	۴/۵۲
۳۶		۴۴۶/۲	۱/۴۹	۶	۱	۷	۳/۰۱	۴/۷۷

ادامه جدول ۳-۳

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
NC1		۴۰۴/۸۷	۲/۰۵	۵	۳	۴	۳/۱	۸/۴۵
NC2		۴۰۶/۴۹	۰/۴۵	۶	۵	۵	۳/۴۷	۷/۶۳
NC3		۴۵۴/۵۵	۰/۸۵	۷	۳	۶	۳/۹۸	۷/۵



ادامه جدول ۳-۳

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pEC <sub>50</sub>
NC4		۴۶۴/۵۴	۰/۴۵	۷	۴	۵	۴/۳۴	۷/۲۳
NC5		۴۵۹/۹۵	۱/۳۵	۶	۲	۶	۴/۱۷	۷/۱۱
NC6		۴۵۳/۵۶	۱/۱۶	۷	۳	۶	۳/۵۴	۶/۹۴
NC7		۴۳۷/۵۶	۲/۴	۶	۳	۳	۳/۸۹	۶/۶۳

علاوه بر بررسی برهم کنش‌های متفاوت هیدروژنی و هیدروفوبی با استفاده از داکینگ مولکولی، پارامترهای قاعده لیپینسکی نیز با استفاده از ابزار وب رایگان Swiss-ADME مورد ارزیابی قرار گرفت. همان‌طور که نتایج جدول ۳-۳ نشان می‌دهد همه ترکیبات پیشنهادی دارای مقادیر قابل قبول می‌باشند و از این نظر خواص فارماکوکینتیکی نیز دارای معیارهای مناسب هستند. پارامتر Syn-Acc نیز نشان می‌دهد که سنتز این ترکیبات در آزمایشگاه امکان‌پذیر است و از درجه آسانی مناسبی برخوردار است.

۳-۱-۳-۴ بررسی میزان و چگونگی تأثیر توصیف‌کننده‌ها بر فعالیت دارویی ترکیبات مجموعه

داده‌های ضد سرطان کارسینوم کولورکتال و ریه و کاربرد مدل LAD-LASSO-LM-ANN

ANN ارائه شده در پیشنهاد ترکیبات جدید

-اثر توصیف‌کننده‌ها بر فعالیت دارویی با استفاده از مدل LAD-LASSO-LM-ANN برای

بازدارنده‌های سرطان کارسینوم کولورکتال

با توجه به نتایج جدول ۲-۱۷، مدل LAD-LASSO-LM-ANN برای بازدارنده‌های سرطان کارسینوم کولورکتال، با ۱۴ توصیف‌کننده در مدل بهینه مشارکت دارند. در ادامه، به شرح مختصری از برخی از توصیف‌کننده‌های منتخب با تفسیرپذیری بیش‌تر پرداخته خواهد شد. توصیف‌کننده nArCO در طبقه گروه‌های عاملی قرار دارد و مربوط به تعداد کتون‌های آروماتیک می‌باشد. علامت منفی این توصیف‌کننده در مدل LAD-LASSO-LM-ANN (جدول ۲-۱۶) نشان می‌دهد که حضور این توصیف‌کننده سبب کاهش مقدار فعالیت دارویی می‌شود. با توجه به ساختارهای مورد مطالعه در مجموعه بازدارنده‌های سرطان کارسینوم کولورکتال، دیده می‌شود که گروه nArCO در ترکیبات شماره ۲۶ تا ۳۳ (فعالیت دارویی ضعیف‌تر با  $pIC_{50}$  کوچک‌تر از ۵/۰) وجود دارد که این موضوع نشان‌دهنده اثر منفی وجود گروه‌های کتون آروماتیک بر فعالیت دارویی بوده که به‌وسیله توصیف‌کننده nArCO وارد مدل شده است. توصیف‌کننده

nArX (X گروه هالوژن است) از طبقه شمارش گروه‌های عاملی می‌باشد. این توصیف کننده تعداد گروه‌های عاملی هالوژنه متصل به حلقه آروماتیک را نشان می‌دهد. ضریب تأثیر مثبت این توصیف کننده نشان می‌دهد که هرچه تعداد گروه‌های هالوژنه متصل به حلقه آروماتیک بیشتر باشد، فعالیت دارویی ترکیب بیشتر می‌شود. توصیف کننده AMW از دسته توصیف‌کننده‌های ساختاری می‌باشد. این توصیف کننده نشان‌دهنده وزن مولکولی متوسط (وزن مولکولی/تعداد کل اتم‌ها) است. این توصیف کننده مستقل از اتصال و ساختار مولکولی بوده و اطلاعاتی از ترکیب مولکولی یک ترکیب را منعکس می‌کند [۲۱۶]. با توجه به ضریب اثر مثبت این توصیف کننده (جدول ۲-۱۶)، مقدار بالاتر این توصیف کننده نشان‌دهنده حضور اتم‌هایی با وزن اتمی بالاتر در مولکول است که منجر به افزایش فعالیت دارویی ترکیباتی با چنین شرایطی می‌شوند. با بررسی بیشتر ترکیبات مورد مطالعه مشاهده می‌شود که ترکیبات شماره ۱۲ (فعالیت دارویی برابر با ۷/۰۹)، ۵۱ (فعالیت دارویی برابر با ۷/۰۰) و ۶۶ (فعالیت دارویی برابر با ۶/۹۲) از دسته ترکیباتی هستند که به دلیل داشتن گروه X (۲ گروه Cl در ترکیبات شماره ۱۲ و ۵۱ و یک گروه Cl در ترکیب شماره ۶۶) دارای AMW بالایی نیز می‌باشند و در نتیجه فعالیت دارویی بالاتری نیز دارند. علاوه بر اثر مثبت توصیف کننده AMW بر متغیر وابسته، توصیف‌کننده‌های تابع توزیع شعاعی (RDF) وزن‌دهی شده با جرم (m) در واقع بیان می‌کنند که این اتم‌ها با وزن بالا باید در محل‌های خاصی حضور داشته باشند و وجود اتم‌هایی با وزن بالا به تنهایی باعث افزایش فعالیت دارویی ترکیب نمی‌شود. در واقع این دسته از توصیف‌کننده‌های RDF، چگونگی توزیع شعاعی این اتم‌ها را در فضای شیمیایی ترکیب نشان می‌دهد. همان‌طور که گفته شد، توصیف‌کننده‌های تابع توزیع شعاعی (RDF)، حاوی اطلاعات اساسی در مورد فاصله‌های بین اتمی در طول پیوند، نوع حلقه و غیره می‌باشند. توصیف‌کننده‌های RDF مبتنی بر توزیع فواصل در نمایش هندسی مولکول و تشکیل یک کد تابع توزیع شعاعی است. RDF020m، RDF105m، RDF130m و RDF135m توصیف‌کننده‌های نرمال شده بر اساس جرم اتمی هستند. این توصیف‌کننده‌ها حضور اتم‌ها را

در کره‌های مجازی با قطرهای ۲۰، ۱۰۵، ۱۳۰ و ۱۳۵ آنگستروم نشان می‌دهند. این ۴ توصیف کننده دارای ضرایب تأثیر متفاوتی در مدل LAD-LASSO-LM-ANN هستند که در جدول ۲-۱۶ آورده شده است. توصیف کننده‌هایی با ضرایب مثبت نشان می‌دهد که وجود اتم‌ها با وزن بالا در ناحیه خاصی باعث افزایش فعالیت دارویی می‌شود و توصیف کننده RDF با ضریب منفی نشان می‌دهد که قرار گیری برخی از اتم‌ها با وزن بالا در ناحیه مشخصی از فضای شیمیایی ترکیب باعث کاهش فعالیت دارویی ترکیبات می‌شود.

-اثر توصیف کننده‌ها بر فعالیت دارویی با استفاده از مدل LAD-LASSO-LM-ANN برای

#### بازدارنده‌های سرطان ریه

مدل LAD-LASSO-ANN ارائه شده برای بازدارنده‌های سرطان ریه (نتایج جدول ۲-۱۷)، مدل با ۷ توصیف کننده دارای شرایط بهینه است. اولین توصیف کننده با بیشترین سهم مشارکت (شکل ۳-۱۴) و بیشترین ضریب تأثیر (جدول ۲-۱۶) مربوط به AMW از دسته توصیف کننده‌های ساختاری می‌باشد. همان‌طور که گفته شد، این توصیف کننده نشان‌دهنده وزن مولکولی متوسط (وزن مولکولی/تعداد کل اتم‌ها) است. با توجه به ضریب اثر مثبت این توصیف کننده (جدول ۲-۱۶)، ترکیباتی با AMW بزرگ‌تر دارای فعالیت دارویی بیش‌تر هستند. توصیف کننده‌های RDF110m و RDF130m مربوط به توصیف کننده‌های تابع توزیع شعاعی هستند که بیانگر اطلاعات اساسی در توزیع فواصل در نمایش هندسی مولکول و تشکیل یک کد تابع توزیع شعاعی است. همان‌طور که گفته شد وجود گروه‌های سنگین در ناحیه‌ای مشخص از مرکز ساختار شیمیایی باعث افزایش فعالیت دارویی ترکیب می‌شود. این توصیف کننده‌ها بر اساس جرم اتمی وزن دهی شده‌اند. بنابراین با توجه به ضریب تأثیر مثبت این دو توصیف کننده (جدول ۲-۱۶)، ترکیباتی با وزن اتمی بیش‌تر، دارای فعالیت دارویی بیش‌تری نیز هستند. MoR21m و MoR12e از دسته توصیف کننده‌های 3D-MoRSE است. MoR21m توصیف کننده سه بعدی وزن دهی شده بر اساس جرم اتمی است. ضریب مثبت این توصیف کننده نیز گواهی بر این موضوع است که ترکیبات با وزن اتمی بیش‌تر

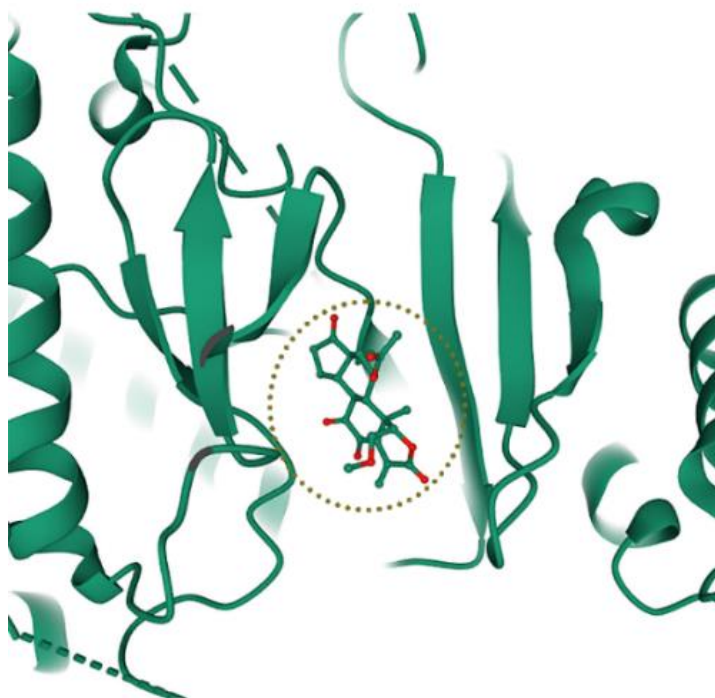
دارای فعالیت دارویی بیش تری می باشند. توصیف کننده Mor12e بر اساس الکترونگاتیوی وزن دهی می شود. Mor12e دارای مقادیر منفی در توصیف کننده و ضریب تأثیر منفی می باشد. بنابراین این توصیف کننده، نشان می دهد که حضور ترکیبات الکترونگاتیو از جمله هالوژن ها در ترکیبات، آن ها را فعال تر می کند [۲۱۷]. بنابراین با بررسی بیش تر ترکیبات مورد مطالعه مشاهده می شود که با توجه به ضریب تأثیر مثبت توصیف کننده های از جنس وزن مولکولی (AMW، RDF110m، RDF130m و MoR21m) و تأثیر مثبت آن ها بر فعالیت دارویی و با توجه به تأثیر مثبت توصیف کننده Mor12e وجود گروه های هالوژنه الکترونگاتیو سنگین هم چون F و Cl در ناحیه مشخص (با توجه به فاصله مکانی مذکور توصیف کننده های RDF و MoRSE) فعالیت دارویی ترکیبات افزایش می یابد. از دسته ترکیبات فعال دارای چند گروه هالوژنه می توان به ترکیبات شماره ۱۲، ۱۳، ۴۹ و ۵۱ اشاره کرد که با داشتن چندین گروه فلوئور و کربن دارای فعالیت دارویی مناسبی نیز هستند (جدول ۲-۱۳ تا جدول ۲-۱۵). nCp توصیف کننده ای از گروه شمارش گروه عاملی می باشد و نشان دهنده کربن  $Sp^3$  انتهایی می باشد. ضریب اثر مثبت این توصیف کننده نشان می دهد که ترکیبات پیشنهادی دارای این گروه کربنی فعالیت دارویی بیش تری را به خود اختصاص داده اند. ترکیبات شماره ۶۸، ۶۷ و ۶۵ با داشتن به ترتیب ۶، ۳ و ۲ گروه متیل از دسته ترکیبات با فعالیت دارویی در محدوده ۶/۶-۶/۹۲ هستند. بنابراین هر چه تعداد گروه متیل روی حلقه فنیل (بال راست) بیش تر باشد فعالیت دارویی آن ترکیب بیش تر است. F084 از گروه قطعات اتم محور است و مربوط به تعداد گروه F متصل به کربن  $Sp^2$  می باشد. با توجه به ضریب اثر مثبت این توصیف کننده، وجود گروه F بر حلقه فنیل (بال راست) ساختار ترکیب باعث افزایش فعالیت دارویی بازدارنده ها می شود. به طور مثال ترکیبات شماره ۱۳ و ۴۹ به ترتیب با داشتن ۱ و ۲ گروه فلوئور بر روی حلقه فنیل از دسته ترکیباتی با فعالیت مناسب دارویی (بزرگ تر از ۶/۰۰) می باشند.

---

<sup>1</sup>Atom centered fragments

## – مطالعه داکینگ مولکولی بازدارنده‌های سرطان کارسینوم کولورکتال و ریه

پس از ایجاد ترکیبات فعال پیشنهادی با فعالیت دارویی مناسب، برهم‌کنش‌های ترکیبات پیشنهادی با اسیدآمینه‌های کلیدی جایگاه فعال گیرنده مورد نظر نیز مورد بررسی قرار گرفت. از این‌رو برهم‌کنش‌های ترکیبات مورد مطالعه و ترکیبات پیشنهادی با استفاده از مطالعه داکینگ مولکولی مورد ارزیابی قرار گرفت. بنابراین، ابتدا ساختار کریستالوگرافی 3HHM به پیشنهاد مقالات منتشر شده از سایت بانک اطلاعاتی پروتئین با فرمت pdb دانلود شد [۸۰، ۸۱، ۱۷۳]. این ساختار دارای ارزش تفکیک  $2/80 \text{ \AA}$  است و به هدف انجام فرایند داکینگ از کیفیت مناسبی برخوردار است [۲۱۸]. به‌منظور استخراج مختصات جایگاه فعال گیرنده، مختصات مرکز ثقل لیگاند کریستالوگرافی به کار گرفته شد. ساختار کریستالوگرافی گیرنده و لیگاند مربوطه در فرایند اعتبار سنجی داکینگ مورد استفاده قرار گرفت. تعداد اجراهای الگوریتم ۱۰۰، ۱۵۰ و ۲۰۰ برای داکینگ لیگاند کریستالوگرافی- گیرنده تعریف شد. پس از انجام فرایند اعتبار سنجی، شرایط بهینه داکینگ، با تعداد اجرای ژنتیک الگوریتم برابر با ۱۵۰ دارای کم‌ترین مقدار RMSD شد. بنابراین در ادامه، داکینگ ترکیبات مورد مطالعه و ترکیبات پیشنهادی در جایگاه فعال گیرنده مورد نظر با ۱۵۰ اجرای ژنتیک الگوریتم داک شدند. فایل خروجی dlg برای همه ترکیبات داک شده به‌دست آمد. با استفاده از اطلاعات مربوط به خروجی داکینگ، بهترین پیکربندی با توجه به کم‌ترین انرژی آزاد اتصال استخراج شدند و برهم‌کنش‌های متفاوت مربوط به لیگاند- گیرنده با استفاده از نرم‌افزار Discovery Studio Visualizer مورد بررسی قرار گرفت.

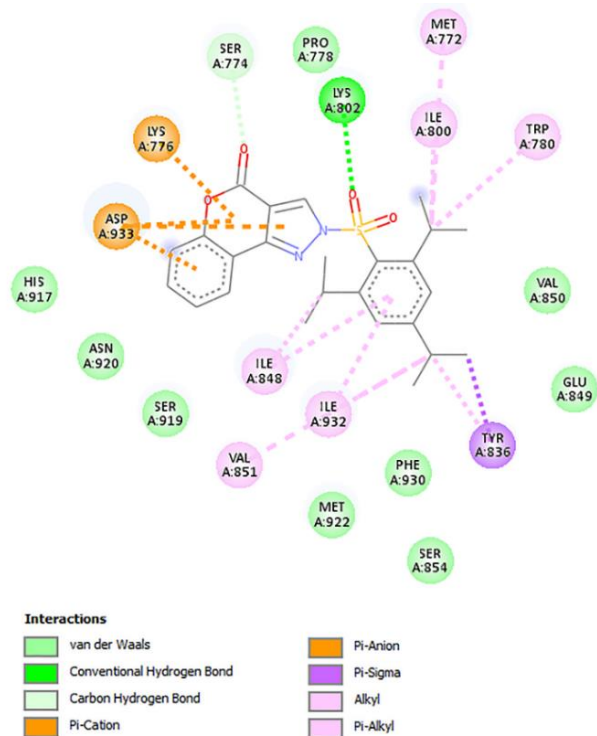


شکل ۳-۲۶ ساختار کریستالوگرافی 3HHM [۲۱۸] (منطقه نقطه چین نشان‌دهنده لیگاند کریستالوگرافی و مابقی زنجیره‌های اسید آمینه‌ای است)

برهم‌کنش‌های متفاوت مربوط به لیگاندهای فعال و کم‌فعال بازدارنده‌های سرطان با گیرنده مربوطه در شکل ۳-۲۷ و شکل ۳-۲۸ آورده شده است. همان‌طور که در شکل ۳-۲۷ مشاهده می‌شود، ترکیب فعال برهم‌کنش‌های متفاوت هیدروژنی و هیدروفوبی با اسید آمینه‌های کلیدی برقرار نموده است. پیوند هیدروژنی ترکیب فعال با Lys802 به‌درستی گواهی بر فعالیت مناسب ترکیب در جایگاه فعال گیرنده می‌باشد. ترکیب کم‌فعال مورد مطالعه توانایی تشکیل پیوندهای هیدروژنی را از خود نشان نداده است. با توجه به نتایج برهم‌کنش‌های این دو ترکیب، فعالیت ترکیبات پیشنهادی نیز مورد قیاس قرار گرفت. به‌طوری‌که نتایج برهم‌کنش ترکیبات پیشنهادی با گیرنده در شکل ۳-۲۹ و شکل ۳-۳۰ نشان داده شده است. با توجه به نتایج مشاهده می‌شود که ترکیب پیشنهادی NC1 علاوه بر این که فعالیت پیش‌بینی شده مناسبی دارد (جدول ۳-۴)، بلکه برهم‌کنش‌های هیدروژنی مناسبی را با اسید آمینه‌های کلیدی از جمله Val851، Glu849 و Gln859 برقرار کرده است. ترکیب پیشنهادی NC2 نیز به‌خوبی توانسته است علاوه

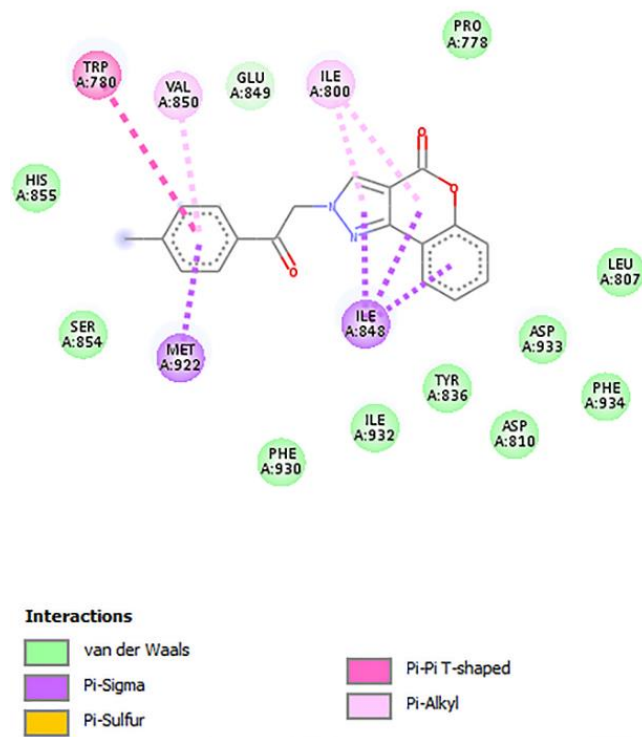
بر پیوندهای هیدروفوبی، با اسید آمینه‌های کلیدی Val851 و Gln859 پیوندهای هیدروژنی برقرار کند. ترکیبات پیشنهادی NC3 و NC4 نیز علاوه بر پیوندهای هیدروفوبی متفاوت، با اسید آمینه Glu859 پیوند هیدروژنی برقرار کنند.

همان‌طور که نتایج نشان داد، ترکیبات پیشنهادی از ارتباط هیدروژنی و هیدروفوبی مناسبی با گیرنده برخوردار هستند و پایداری این ترکیبات در جایگاه فعال گیرنده مورد نظر، مورد تأیید است. علاوه بر بررسی برهم کنش‌های ترکیبات، ویژگی‌های فارماکوکینتیکی ترکیبات نیز مورد محاسبه قرار گرفت و نتایج جدول ۳-۴ نشان می‌دهد که ترکیبات پیشنهادی از نظر قاعده لیپینسکی نیز مورد تأیید هستند و علاوه بر این درجه آسانی سنتز آزمایشگاهی کم‌تر از ۱۰ این ترکیبات بالقوه نیز محاسبه شد و با توجه به نتایج جدول ۳-۴ مشخص است که سنتز این ترکیبات با پیچیدگی همراه نیست.

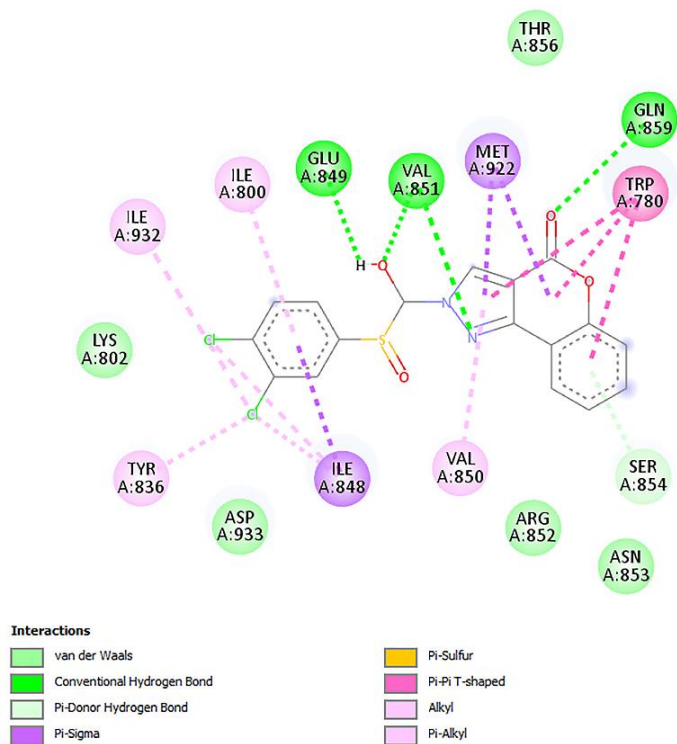


شکل ۳-۲۷ برهم‌کنش ترکیب نسبتاً فعال (۶۸) موجود در مجموعه داده‌های ضد سرطان کاسینوم کولورکتال با اسید آمینه‌های کلیدی

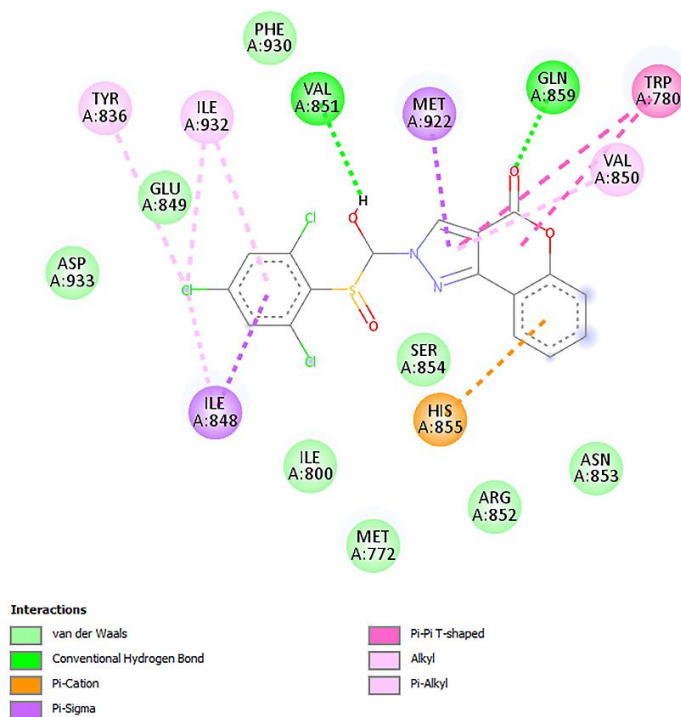




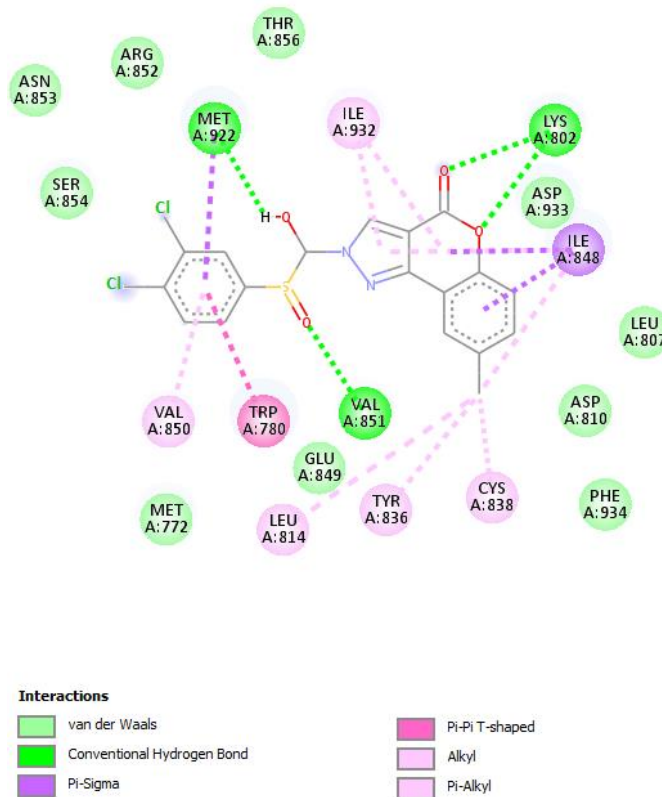
شکل ۳-۲۸ برهم کنش ترکیب کم فعال (۲۷) موجود در مجموعه داده‌های ضد سرطان کاسینوم کولورکتال با اسید آمینه‌های کلیدی



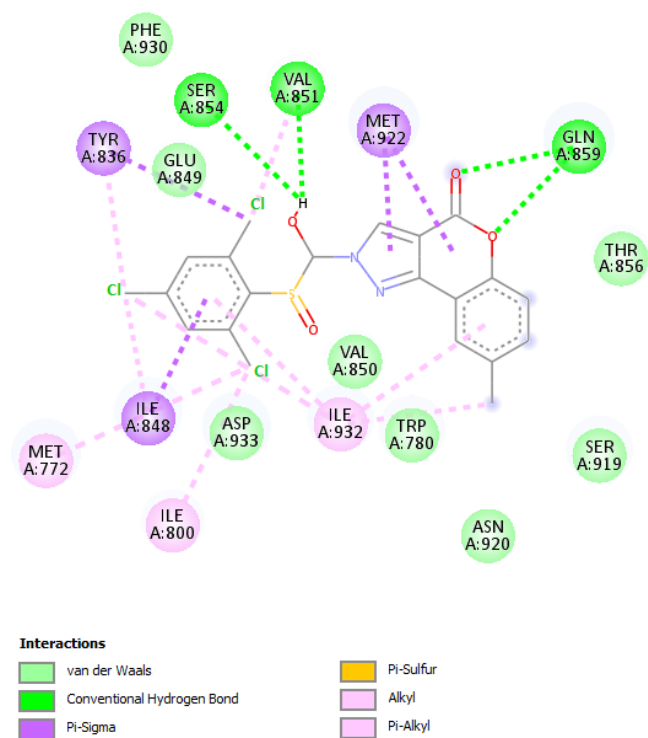
شکل ۳-۲۹ برهم کنش ترکیب پیشنهادی NC1 با اسید آمینه‌های کلیدی



شکل ۳-۳ برهم کنش ترکیب پیشنهادی NC2 با اسید آمینه‌های کلیدی

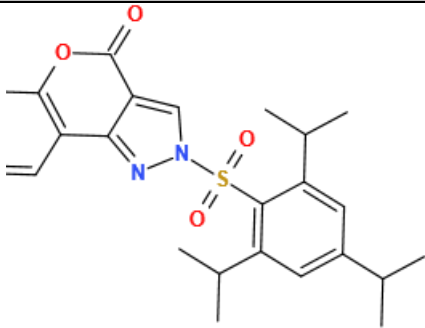
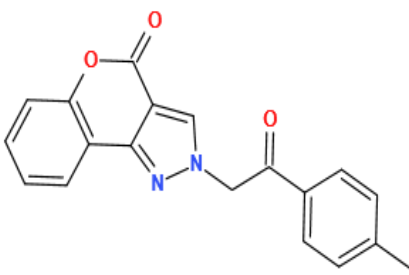
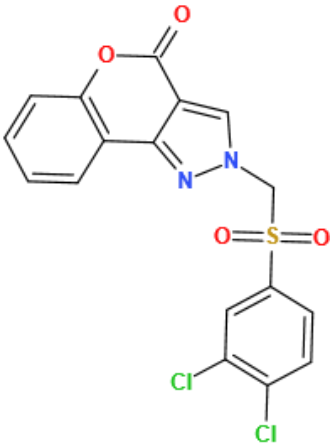


شکل ۳-۳ برهم کنش ترکیب پیشنهادی NC3 با اسید آمینه‌های کلیدی

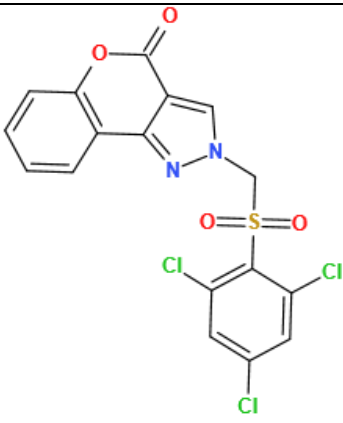
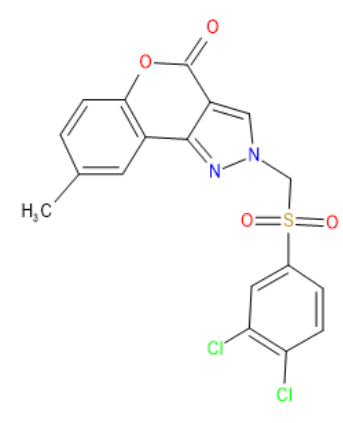
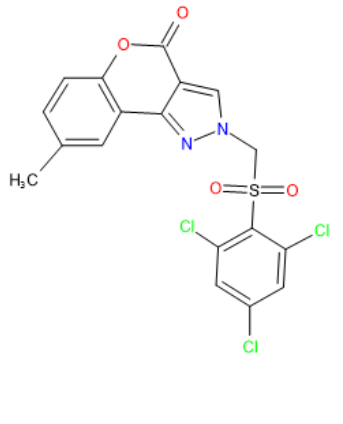


شکل ۳-۳۲ برهم‌کنش ترکیب پیشنهادی NC4 با اسید آمینه‌های کلیدی

جدول ۳-۴ پارامترهای PK محاسبه شده برای ترکیبات مورد مطالعه ضد سرطان و ترکیبات پیشنهادی

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pIC <sub>50</sub>
۶۸		۳۱۸/۳۳	۲/۶۲	۳	۰	۴	۲/۸۸	۷/۱۵
۲۶		۴۵۲/۵۷	۵/۲	۵	۰	۵	۴/۶۸	۴/۲۰
NC1		۴۰۹/۲۴	۳/۳۵	۳	۰	۵	۳/۵۱	۸/۹۲

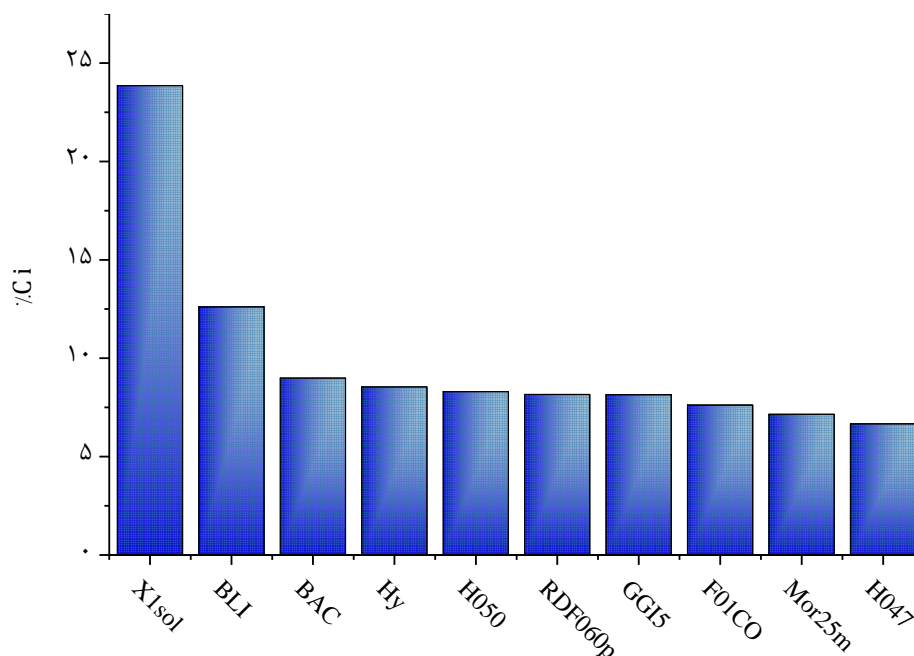
ادامه جدول ۴-۳

شماره ترکیب	ساختار شیمیایی	MW	MLOGP	#Rot-B	#H-B-don	#H-B-acc	Syn-Acc	pIC <sub>50</sub>
NC2		۴۴۳/۶۹	۳/۸۵	۳	۰	۵	۳/۵۸	۸/۷۸
NC3		۴۵۷/۷۰	۴/۰۳	۳	۰	۵	۳/۵۸	۷/۵۲
NC4		۴۲۳/۲۰	۴/۱۳	۳	۰	۵	۳/۵۱	۶/۴۳

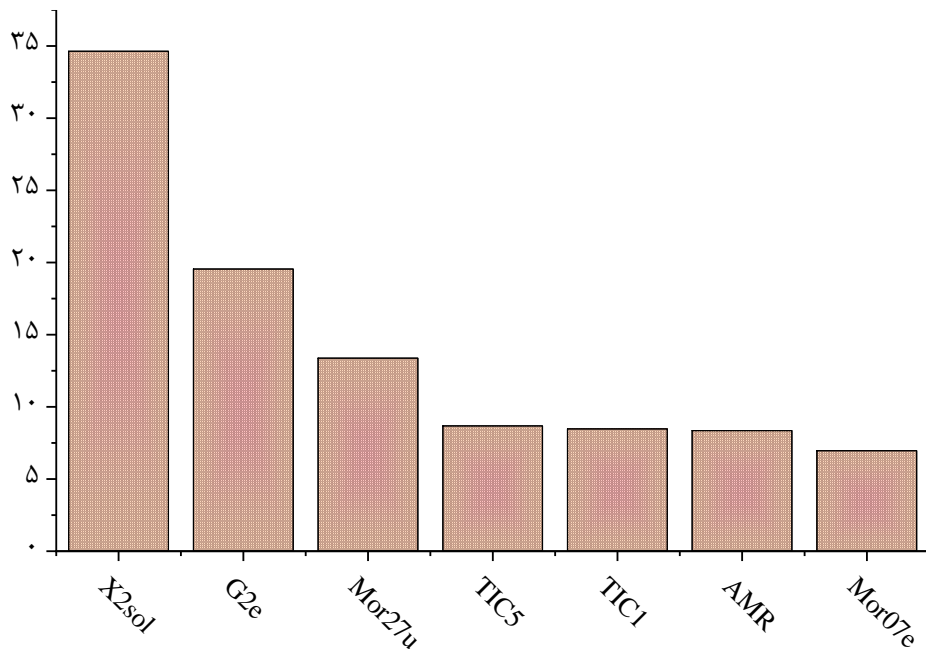
### ۳-۱-۴ تحلیل توصیف‌کننده‌های مدل SCAD-ANN برای ترکیبات آلی فرار

#### ۳-۱-۴-۱ محاسبه سهم مشارکت هر توصیف‌کننده در مدل SCAD-ANN

سهم مشارکت توصیف‌کننده‌های منتخب SCAD در مدل غیرخطی توسعه‌یافته بهینه‌ی مورد بررسی قرار گرفت. برای محاسبه میزان درصد مشارکت هر توصیف‌کننده، مقادیر هر توصیف‌کننده در محدوده تغییرات مقادیر واقعی هر توصیف‌کننده تصادفی شد. این بار مدل‌های بهینه SCAD-ANN هر بار در حضور توصیف‌کننده  $i$  با مقادیر تصادفی و سایر توصیف‌کننده‌ها با مقادیر واقعی خود توسعه داده شدند و مقادیر RI با این شرایط پیش‌بینی شد. مقادیر  $RMSE_i$  مربوط به مجموعه ارزیابی در حضور توصیف‌کننده  $i$  (با مقادیر تصادفی) به دست آمد و این فرایند برای همه توصیف‌کننده‌ها تکرار شد و در نهایت درصد مشارکت هر توصیف‌کننده ( $\%C_i$ ) بر اساس رابطه ۱-۱۵ محاسبه شد و نتایج در شکل ۳-۳۳ و شکل ۳-۳۴ نمایش داده شد در ادامه به بررسی میزان و چگونگی تأثیر این توصیف‌کننده‌ها بر شاخص بازدارندگی ترکیبات و رابطه آن‌ها با ساختار ترکیبات مورد مطالعه پرداخته خواهد شد.



توصیف‌کننده‌های منتخب روش SCAD برای مجموعه داده‌های A  
شکل ۳-۳۳ نمودار سهم مشارکت توصیف‌کننده‌ها در مدل SCAD-ANN برای مجموعه A



توصیف کننده های منتخب روش SCAD برای مجموعه داده های B  
 شکل ۳-۳ نمودار سهم مشارکت توصیف کننده ها در مدل SCAD-ANN برای مجموعه B

### ۳-۱-۴-۲ بررسی رابطه بین توصیف کننده های استفاده شده در مدل SCAD-ANN و شاخص

#### بازداری ترکیبات مورد مطالعه

برای بررسی بیشتر چگونگی تأثیر هر توصیف کننده بر شاخص بازداری، علامت تأثیر هر توصیف کننده بر RI تعیین شد. برای این کار، ضرایب مدل SCAD مطابق با معادله های زیر برای مجموعه داده های A و B (رابطه ۳-۳) و B (رابطه ۴-۳) به دست آمد.

$$RI = 0.35 + 0.97X_{1sol} - 0.002F_{01CO} - 0.22BAC - 0.01H_{047} - 0.03Mor_{25m} + 0.03H_{050} + 0.03Hy$$

رابطه ۳-۳

$$- 0.25BLI - 0.02RDF_{060p} - 0.01GGI_5 + 0.015F_{03CN} + 0.016C_{025} + 0.02HATS_{5V} + 0.03TPSAT_{Tot} + 0.001E_{1m}$$

$$RI = 0.42 + 1.01X_{2sol} - 0.005G_{2e} - 0.13Mor_{27u} + 0.027TIC_5 + 0.006TIC_1 + 0.005AMR$$

رابطه ۴-۳

$$- 0.005Mor_{07e}$$

با توجه به رابطه ۳-۳ و رابطه ۴-۳، اثر افزایشی و یا کاهشی تأثیر توصیف کننده ها بر شاخص بازداری نشان داده شد. با توجه به تفسیرپذیری بهتر و تأثیر بیش تر برخی از توصیف کننده ها به شرح آنها

پرداخته شد و در ادامه شرحی از این توصیف کننده‌ها و چگونگی وابستگی شاخص بازداری به آن‌ها ارائه شده است.

#### -اثر توصیف کننده‌ها بر شاخص بازداری با استفاده از مدل SCAD برای مجموعه A

با توجه به نمودار سهم مشارکت توصیف کننده‌ها مؤثرترین توصیف کننده‌ها عبارتند از X1sol، BAC و BLI و این توصیف کننده‌ها به ترتیب اثر افزایشی، کاهش و کاهش بر RI دارند. بحرانی‌ترین توصیف کننده با بیش‌ترین سهم در مدل پیشنهادی مربوط به X1sol است. این توصیف کننده از دسته توصیف کننده‌های مکانی است و شاخص اتصال حلالیت 1-chi نام دارد [۲۱۹]. توصیف کننده‌های شاخص اتصال انحلال پذیری، بیانی از آنتروپی هستند. این توصیف کننده بیانی از پارامترهای اتمی است و با رابطه ۳-۵ نشان داده می‌شود:

$$m_{\chi^s} = \left(\frac{1}{2^{m+1}}\right) \cdot \sum_{k=1}^k \left( \frac{(\prod_{a=1}^k L_a)_k}{(\prod_{a=1}^n \delta_a)^{\frac{1}{2}k}} \right) \quad \text{رابطه ۳-۵}$$

که در آن  $L_a$  عدد کوانتومی اصلی (۲ برای اتم‌های C، N، O و ۳ برای Si، S، Cl، ...) مربوط به  $a_{th}$  اتم در زیرگراف  $k_{th}$  است.  $\delta_a$  درجه رأس مربوطه است.  $k$  تعداد کل زیرگراف‌های مرتبه  $m$  ام و  $n$  تعداد رئوس در زیرگراف است.  $1/2^{m+1}$  ضریب نرمال سازی است و  $m$  نشان دهنده درجه شاخص اتصال حلالیت است که برای X1sol برابر با ۱ است [۹]. با توجه به علامت مثبت ضریب X1sol در معادله SCAD، ترکیبات با مقادیر X1sol زیاد دارای مقادیر RI بالاتری نیز هستند.

BLI توصیف کننده مورد بحث بعدی است و با توجه به رابطه ۳-۵ دارای ضریب اثر منفی بر شاخص بازداری است. شاخص شباهت بنزن کایر (Keir) یک شاخص آروماتیکی است که از توپولوژی مولکولی محاسبه می‌شود. توصیف کننده‌های BLI از جمله شاخص‌های رزونانسی هستند که پایداری بنزن را از نظر تئوری توضیح می‌دهند و میزان جابجایی سیستم‌های مزدوج را پیش‌بینی می‌کنند.



شاخص BAC از دسته توصیف‌کننده‌های مکانی و مربوط به شاخص مرکزی balaban می‌باشد. نام شاخص مرکزی به این دلیل است که توپولوژی نمودار را از مرکز مشاهده می‌کند. این توصیف‌کننده پیچیدگی و انشعاب مولکولی را رمزگذاری می‌کند. پارامتر BAC می‌تواند به یک یا چند ویژگی ساختاری مولکول مانند اندازه، شکل، تقارن، انشعاب و چرخه‌ای بودن حساس باشد. همچنین می‌تواند اطلاعات شیمیایی مربوط به نوع اتم و تعداد پیوند را نیز رمزگذاری کند.

### -اثر توصیف‌کننده‌ها بر شاخص بازداری با استفاده از مدل SCAD برای مجموعه B

از بین توصیف‌کننده‌های مدل بهینه برای مجموعه B، X2sol، G2e و Mor27u بیش‌ترین سهم مشارکت را در ساخت مدل ANN دارند. از این‌رو در ادامه به شرح مفصلی از هر یک از این‌ها پرداخته می‌شود.

بحرانی‌ترین توصیف‌کننده با بیش‌ترین سهم مشارکت در مدل پیشنهادی X2sol است که نشان‌دهنده شاخص اتصال حلالیت با درجه ۲ است. این توصیف‌کننده از دسته توصیف‌کننده‌های توپولوژیکی دوبرعی است [۲۱۹]. توضیحات مربوط به این شاخص در بخش قبلی و با رابطه ۳-۵ شرح داده شد. با توجه به علامت مثبت ضریب X2sol در معادله SCAD (رابطه ۳-۴)، ترکیبات با مقادیر X2sol زیاد دارای مقادیر RI بالاتری هستند.

دومین توصیف‌کننده مهمی (شکل ۳-۳۴) که در مدل ظاهر شده است G2e است که متعلق به شاخص WHIM است. این توصیف‌کننده، مؤلفه‌ای است که توسط الکترونگاتیوی ساندerson اتمی وزن شده است. G2e مربوط به شاخص‌های آماری است و بر روی پیش‌بینی اتم‌ها در امتداد محورهای اصلی محاسبه می‌شود. الکترونگاتیوی ساندerson اتمی یکی از طرح‌های وزن دهی است که برای محاسبه ماتریس کوواریانس وزنی در توصیف‌کننده G2e استفاده می‌شود و به شکل مولکولی و اندازه کل ساختار مربوط می‌شود [۲۲۰]. ترکیب با مقدار G2e کم‌تر، استخلاف بیش‌تری دارد. بنابراین، اندازه ساختار بزرگ‌تر است.

با توجه به تأثیر منفی G2e بر روی RI، ترکیب با مقدار G2e کم‌تر دارای مقدار RI بیش‌تری است. این مفهوم نشان می‌دهد که افزایش تعداد استخلاف‌ها و حجم کل ساختار باعث می‌شود فاز ساکن با ترکیب مورد مطالعه بیش‌تر برهمکنش داشته باشد و مقدار RI افزایش یابد.

سومین توصیف‌کننده (Mor27u) به دسته توصیف‌کننده‌های نمایش ساختارهای سه بعدی مولکولی بر اساس پراش الکترونی (3D-MoRSE) تعلق دارد. پارامترهای عددی 3D-MoRSE ساختار مولکولی سه بعدی را به‌طور قابل توجهی نشان می‌دهند و آرایش سه بعدی اتم‌ها را منعکس می‌کنند و پیوندهای شیمیایی در نظر گرفته نمی‌شوند [۲۲۱]. Mor27u (3D-MoRSE-signal 27/unweighted) مبتنی بر ایده به‌دست آوردن اطلاعات از مختصات اتمی سه بعدی با تبدیل مورد استفاده در مطالعات پراش الکترون برای تهیه منحنی‌های پراکندگی نظری است [۲۲۲]. به‌طور کلی توصیف‌کننده‌های 3D-MoRSE فواصل بین اتمی را اندازه‌گیری می‌کنند و دیدگاه‌های متفاوتی از کل ساختار مولکول را نشان می‌دهند، اگرچه معنای آن‌ها هنوز غیر قابل مشخص است [۲۲۳]. با توجه به تأثیر منفی این توصیف‌کننده بر مقادیر RI، انتظار می‌رود که ترکیبات با مقادیر منفی Mor27u دارای مقادیر RI نسبتاً بالایی باشند. در نتیجه، ظهور Mor27u در مدل نهایی، ویژگی‌های ساختاری مولکول‌های با اندازه بزرگ را در مدل در نظر می‌گیرد.

## ۲-۳ نتیجه‌گیری نهایی

در این مطالعه، مدل‌های QSAR\QSRR مبتنی بر به‌کارگیری ترکیب روش‌های انقباضی و مدل غیر خطی شبکه عصبی مصنوعی برای پیش‌بینی فعالیت دارویی/شاخص بازدارندگی مجموعه داده‌های متفاوتی مورد استفاده قرار گرفت. در ادامه نتیجه‌گیری کلی از نتایج حاصل از هر یک از مدل‌های توسعه یافته QSAR و QSRR ارائه خواهد شد.

در مطالعه اول این رساله، برای اولین بار، روش انقباضی SCAD (فن و لی، ۲۰۰۱) به عنوان روش انتخاب متغیر انقباضی با مدل غیر خطی LM-ANN ترکیب شد و در مطالعات QSAR مورد استفاده قرار گرفت. روش پیشنهادی SCAD-LM-ANN از مزایای ذاتی SCAD در کاهش ابعاد داده‌ها و کارایی LM-ANN در مدل‌سازی روابط غیر خطی بهره برده است. در نتیجه، مدل ساخته شده علاوه بر قدرت پیش‌بینی رضایت بخش از تفسیر پذیری و تنگی مناسبی نیز برخوردار است. مدل توسعه یافته SCAD-LM-ANN با استفاده از ارزیابی نتایج آماری مجموعه آزمون، تکنیک LOO، دامنه کاربرد، آزمون Y-تصادفی و محاسبه پارامترهای آماری مورد ارزیابی بیش‌تر قرار گرفت. همان‌طور که نتایج (جدول ۲-۶) نشان می‌دهد، مدل از قدرت پیش‌بینی و تعمیم‌پذیری قابل قبولی برخوردار است. توانایی ذاتی SCAD در انتخاب متغیر منجر به ساخت یک مدل تفسیرپذیر با مجموعه کوچکی از توصیف‌کننده‌ها می‌شود. توصیف‌کننده‌های مناسب و صحیح انتخاب شده توسط روش SCAD، امکان طراحی ترکیبات فعال جدید را با توجه به تأثیر مقادیر ضرایب رگرسیونی توصیف‌کننده‌ها بر فعالیت دارویی فراهم می‌کند. صحت فعال بودن ترکیبات پیشنهادی، با استفاده از بررسی برهم‌کنش لیگاند-گیرنده از طریق مطالعه داکینگ مولکولی تأیید شد. مقادیر پیش‌بینی شده فعالیت دارویی (جدول ۳-۱)، نتایج برهم‌کنش داکینگ مولکولی ترکیبات پیشنهادی با گیرنده (شکل ۳-۳ و شکل ۴-۳) و نتایج ویژگی‌های فارموکینتیکی ترکیبات پیشنهادی (جدول ۳-۱) نشان می‌دهد که ترکیبات جدید فعالیت بالقوه‌ای را به‌عنوان بازدارنده‌های جدید ایدز دارند [۱۱۴].

در مطالعه دوم این رساله، با توجه به اهمیت بررسی مهارکننده‌های 3CL<sup>pro</sup> برای مقابله با بیماری کووید-۱۹، مدل شبکه عصبی با روش انتخاب متغیر انقباضی قدرتمند ALASSO (زو، ۲۰۰۶)، برای پیش‌بینی فعالیت بازدارنده‌های ضد کووید-۱۹، توسعه داده شد. روش انقباضی ALASSO تعداد توصیف کننده‌ها را به ۹ توصیف کننده مؤثر با بیش‌ترین تأثیر بر متغیر وابسته کاهش داد. مدل ALASSO-LM-ANN با قابلیت تفسیرپذیری و قدرت پیش‌بینی مناسب تولید شد. ارزیابی مدل ALASSO-LM-ANN با استفاده از روش‌های مختلف برآورد شد. نتایج خلاصه شده در جدول ۲-۱۲ نشان می‌دهد که مدل ALASSO-LM-ANN دارای پارامترهای آماری قابل قبولی است. همچنین، نمودار دامنه کاربرد (شکل ۲-۱۴) به وضوح اعتبار مدل توسعه یافته ALASSO-LM-ANN را اثبات می‌کند. بنابراین، کارایی مدل تفسیرپذیر و پیشگو توسعه یافته با توصیف کننده‌های منتخب روش قدرتمند ALASSO برای پیشنهاد ترکیبات جدید مورد استفاده قرار گرفت. صحت ترکیبات پیشنهادی با استفاده از بررسی بر هم کنش ترکیبات پیشنهادی - گیرنده به‌کمک داکینگ مولکولی و محاسبه خواص PK برای تأیید قانون پنج لیپینسکی مورد ارزیابی قرار گرفت. مقایسه ترکیبات فعال مورد مطالعه و ترکیبات پیشنهادی به درستی نشان دهنده قدرت بالای ترکیبات پیشنهادی به‌عنوان بازدارنده‌های بالقوه کووید-۱۹ هستند [۱۱۵].

در مطالعه سوم این رساله، این مدل‌های LAD-LASSO-LM-ANN برای پیش‌بینی فعالیت دارویی سه مجموعه داده متفاوت ضد ایدز و ضد سرطان مورد استفاده قرار گرفت. موثرترین توصیف کننده‌های هر مجموعه داده با استفاده از روش LAD-LASSO (وانگ، ۲۰۰۷) به‌عنوان ورودی ANN تعریف شدند. پس از ساخت و توسعه مدل‌های ANN، فعالیت‌های دارویی ترکیبات مجموعه آزمون بر اساس مدل بهینه LAD-LASSO-LM-ANN پیش‌بینی شدند. پارامترهای آماری متفاوتی برای مدل‌های توسعه یافته محاسبه شد (جدول ۲-۲۲) و نتایج همه پارامترهای آماری در محدوده قابل قبول قرار داشتند. در این مطالعه، چگونگی وابستگی فعالیت دارویی به توصیف کننده‌های منتخب هر مجموعه داده مورد بررسی قرار

گرفت. با توجه به ارتباط بین توصیف کننده‌های منتخب و فعالیت دارویی ترکیبات، چندین بازدارنده جدید با فعالیت دارویی مناسب پیشنهاد شد. صحت ترکیبات پیشنهادی با استفاده از بررسی برهم کنش ترکیبات پیشنهادی - گیرنده (شکل ۳-۱۶ تا شکل ۳-۳۲) و محاسبه خواص PK برای تأیید قانون پنج لیپینسکی مورد ارزیابی قرار گرفت (جدول ۳-۳ و جدول ۳-۴). مقایسه ترکیبات فعال مورد مطالعه و ترکیبات پیشنهادی به درستی نشان دهنده فعالیت دارویی مناسب ترکیبات پیشنهادی است. علاوه بر این، فاکتور سهولت سنتز (کمتر از ۱۰) نشان می‌دهد که سنتز ترکیبات در مقیاس آزمایشگاهی نیز امکان پذیر است. در مطالعه آخر این رساله تحقیقاتی، از روش SCAD به عنوان یک روش انقباضی کارآمد برای انتخاب توصیف کننده‌های مؤثر و پیش بینی مقادیر RI ترکیبات آلی فرار متفاوت با استفاده از مدل غیر خطی ANN برای ساخت مدل QSRR استفاده شد. مدل‌های متناظر SCAD-ANN پیشنهادی به طور هم‌زمان از مزیت SCAD مانند تنگی و پایداری مناسب توام با توانایی پیش‌بینی بسیار بالای روش شبکه عصبی مصنوعی استفاده شد. ارزیابی مدل‌های SCAD-ANN با استفاده از داده‌های مجموعه آزمون، تکنیک LOO، دامنه کاربرد و آزمون Y - تصادفی انجام شد و نتایج نشان دهنده قدرت پیش‌بینی و تعمیم پذیری رضایت‌بخش مدل ارائه شده است.

در نهایت، به منظور بررسی کارایی روش‌های انقباضی به‌عنوان روش‌های مدل‌سازی QSAR و QSPR، هر کدام از روش‌های مورد مطالعه در شرایط یکسان و مشابه با مراحل ارائه شده در هر بخش، بر روی داده‌های آموزش مجموعه داده‌های متفاوت اعمال شد. فعالیت‌های دارویی / شاخص‌های بازدارندگی ترکیبات آزمون با استفاده از مدل‌های توسعه یافته با توصیف کننده‌های مولکولی منتخب پیش‌بینی شد. پارامترهای آماری متفاوت (جدول ۱-۱) برای ترکیبات مجموعه آزمون پیش‌بینی شده، محاسبه شد. با مقایسه نتایج مربوط به به‌کارگیری روش‌های انقباضی جفت شده با مدل شبکه عصبی مصنوعی و روش‌های انقباضی به‌عنوان روش‌های انتخاب متغیر و مدل‌سازی هم‌زمان، مشاهده می‌شود که مدل‌های غیرخطی

ANN جفت شده با تکنیک‌های جریمه‌ای (SCAD, ALASSO و LAD-LASSO) از قدرت پیش‌بینی و تعمیم‌پذیری بهتری برخوردار هستند و در برخی از موارد به کارگیری مدل‌های انقباضی به‌تنهایی نتوانسته است نتایج قابل‌قبولی را ارائه دهد. بنابراین هدف این رساله در توسعه مدل‌های غیرخطی ANN جفت شده با تکنیک‌های جریمه‌ای به‌عنوان روش‌های انتخاب‌متغیر، به‌طور کارآمد و موفق عمل نموده است.

در نهایت می‌توان به این نتیجه رسید که برای ساخت مدل‌های QSAR/QSPR با قدرت پیش‌بینی بالا، مدل SCAD و LAD-LASSO عملکرد بالاتری را از خود نشان داده‌اند و نتایج رضایت‌بخشی دارند. در حالی که برای ساخت مدل‌های QSAR/QSPR تنک و تفسیرپذیر روش ALASSO برای انتخاب موثرترین توصیف‌کننده‌ها پیشنهاد می‌شود ولی از آنجایی که ALASSO به دنبال برآورد ضرایب حقیقی و درست‌تری است، یافتن مدل QSAR/QSPR با بیش‌ترین ارتباط بین متغیرهای مستقل و وابسته کمی دشوار است.

جدول ۳-۵ مقایسه پارامترهای آماری محاسبه شده برای مجموعه آزمون مدل برتر SCAD-LM-ANN با مدل SCAD برای مشتقات استانیلید/ استامید به عنوان بازدارنده‌های ایدز

ردیف	پارامترهای آماری	مقادیر $pEC_{50}$ پیش بینی شده	مقادیر $pEC_{50}$ پیش بینی شده
		ترکیبات مجموعه آزمون با مدل SCAD	ترکیبات مجموعه آزمون با مدل SCAD-LM-ANN
۱	PRESS	۲/۶۳	۱/۳۶
۲	SEP	۰/۴۹	۰/۳۵
۳	MAE	۰/۴۲	۰/۳۱
۴	REP(%)	۸/۰۵	۵/۷۹
۵	MSE	۰/۲۴	۰/۱۲
۶	MRE	۷/۰۷	۵/۱
۷	$R^2$	۰/۶۵	۰/۹۲
۸	$R_0^2$	۰/۴۸	۰/۹۱
۹	$R_0^2$ نسبی	۰/۲۶	۰/۰۱
۱۰	$R_m^2$	۰/۳۸	۰/۹۲
۱۱	$R_0'^2$	۰/۶۵	۰/۸۹
۱۲	$R_0'^2$ نسبی	۰/۰۰۲	۰/۰۳
۱۳	$R_m'^2$	۰/۶۳	۰/۹۲
۱۴	R-R	۰/۱۷	۰/۰۲
۱۵	k	۱/۰۰	۰/۹۶
۱۶	k'	۰/۹۹	۱/۰۴

جدول ۳-۶ مقایسه پارامترهای آماری محاسبه شده برای مجموعه آزمون مدل برتر ALASSO-LM-ANN با مدل ALASSO برای مشتقات 3-chymotrypsin like protease (3CLPro) به‌عنوان بازدارنده‌های SARS-COV-2

ردیف	پارامترهای آماری	مقادیر $pIC_{50}$ پیش بینی شده ترکیبات مجموعه آزمون با مدل ALASSO	مقادیر $pIC_{50}$ پیش بینی شده ترکیبات مجموعه آزمون با مدل ALASSO-LM-ANN
۱	PRESS	۱۴/۱۱	۱/۴۷
۲	SEP	۰/۴۰	۰/۳۴
۳	MAE	۰/۳۰	۰/۲۸
۴	REP(%)	۷/۶۹	۶/۰۸
۵	MSE	۰/۱۶	۰/۱۱
۶	MRE	۵/۸۳	۵/۱۴
۷	$R^2$	۰/۷۵	۰/۸۳
۸	$R_0^2$	۰/۵۸	۰/۸۵
۹	$R_0^2$ نسبی	۰/۲۳	۰/۰۱
۱۰	$R_m^2$	۰/۴۴	۰/۷۷
۱۱	$R_0'^2$	۰/۷۴	۰/۸۶
۱۲	$R_0'^2$ نسبی	۰/۰۱	۰/۰۰
۱۳	$R_m'^2$	۰/۳۵	۰/۷۷
۱۴	R-R	۰/۱۶	۰/۰۱
۱۵	k	۰/۹۹	۰/۹۸
۱۶	k'	۱/۰۰	۱/۰۳



جدول ۳-۷ مقایسه پارامترهای آماری محاسبه شده برای مجموعه آزمون مدل برتر LAD-LASSO-ANN با مدل LAD- LASSO برای هر سه مجموعه داده‌های متفاوت

ردیف	پارامتر آماری	مجموعه داده‌های ضد سرطان			مجموعه داده‌های ضد سرطان ریه		
		LAD- LASSO	LAD- LASSO- ANN	LAD- LASSO	LAD- LASSO- ANN	LAD- LASSO	LAD- LASSO- ANN
۱	PRESS	۱/۷۸	۱/۴۸	۱/۳۷	۰/۸۰	۴/۲۱	۱/۱۷
۲	SEP	۰/۴۰	۰/۳۷	۰/۳۵	۰/۲۷	۰/۶۲	۰/۳۳
۳	MAE	۰/۳۱	۰/۲۹	۰/۲۷	۰/۲۳	۰/۴۵	۰/۲۸
۴	REP(%)	۵/۷۸	۵/۶۲	۶/۷۶	۵/۱۴	۱۲/۴۳	۶/۵۴
۵	MSE	۰/۱۶	۰/۱۳	۰/۱۲	۰/۰۷	۰/۳۸	۰/۱۱
۶	MRE	۴/۵۶	۵/۲۴	۴/۹۰	۴/۴۲	۹/۱۷	۵/۸۹
۷	R <sup>2</sup>	۰/۸۲	۰/۸۷	۰/۶۷	۰/۸۴	۰/۵۲	۰/۸۷
۸	R <sub>0</sub> <sup>2</sup>	۰/۷۹	۰/۸۶	۰/۵۹	۰/۷۷	۰/۴۲	۰/۸۵
۹	R <sub>0</sub> <sup>2</sup> نسبی	۰/۰۴	۰/۰۱	۰/۱۲	۰/۰۸	۰/۱۹	۰/۰۲
۱۰	R <sub>m</sub> <sup>2</sup>	۰/۶۸	۰/۷۸	۰/۴۸	۰/۶۲	۰/۳۶	۰/۷۵
۱۱	R <sub>0</sub> ' <sup>2</sup>	۰/۸۰	۰/۸۰	۰/۶۶	۰/۸۴	۰/۵۱	۰/۸۷
۱۲	R <sub>0</sub> ' <sup>2</sup> نسبی	۰/۰۲	۰/۰۸	۰/۰۱	۰/۰۰	۰/۰۲	۰/۰۰
۱۳	R <sub>m</sub> ' <sup>2</sup>	۰/۷۱	۰/۶۵	۰/۴۸	۰/۵۷	۰/۲۹	۰/۷۳
۱۴	R-R	۰/۰۱	۰/۰۶	۰/۰۷	۰/۰۷	۰/۰۹	۰/۰۲
۱۵	k	۰/۹۷	۰/۹۷	۰/۹۹	۰/۹۸	۰/۹۸	۰/۹۸
۱۶	k'	۱/۰۳	۱/۰۳	۱/۰۰	۱/۰۲	۱/۰۰	۱/۰۱

جدول ۳-۸ مقایسه پارامترهای آماری محاسبه شده برای مجموعه آزمون مدل برتر SCAD-ANN با مدل SCAD برای مجموعه داده A و B

ردیف	پارامترهای آماری	مجموعه داده A		مجموعه داده B	
		SCAD	SCAD-ANN	SCAD	SCAD-ANN
۱	PRESS	۰/۱۰	۰/۰۵	۲۸۴۰۶۰	۷۸۴۲۶
۲	SEP	۰/۰۶	۰/۰۴	۱۶۰/۷۰	۸۴/۴۴
۳	MAE	۰/۰۴	۰/۰۳	۱۲۶/۸۴	۶۶/۹۱
۴	REP(%)	۱۰/۴۱	۵/۳۹	۱۵/۲۸	۸/۰۲
۵	MSE	۰/۰۶	۰/۰۴	۱۶۰/۷۰	۸۴/۴۳
۶	MRE	۶/۷۴	۴/۳۴	۱۲/۱۳	۶/۵۳
۷	R <sup>2</sup>	۰/۸۳	۰/۹۲	۰/۷۳	۰/۸۹
۸	R <sub>0</sub> <sup>2</sup>	۰/۸۲	۰/۹۱	۰/۷۳	۰/۸۹
۹	R <sub>0</sub> <sup>2</sup> نسبی	۰/۰۱	۰/۰۱	۰/۰۰	۰/۰۱
۱۰	R <sub>m</sub> <sup>2</sup>	۰/۷۵	۰/۸۳	۰/۷۳	۰/۸
۱۱	R <sub>0</sub> <sup>'2</sup>	۰/۸۱	۰/۹۲	۰/۶۸	۰/۸۸
۱۲	R <sub>0</sub> <sup>'2</sup> نسبی	۰/۰۲	۰/۰۰	۰/۰۷	۰/۰۰
۱۳	R <sub>m</sub> <sup>'2</sup>	۰/۷۴	۰/۸۲	۰/۵۷	۰/۷۹
۱۴	R-R	۰/۰۱	۰/۰۱	۰/۰۵	۰/۰۱
۱۵	k	۱/۰۰	۰/۹۸	۱/۰۸	۱/۰۴
۱۶	k'	۰/۹۸	۱/۰۱	۰/۹۱	۰/۹۵

### ۳-۳ آینده نگری

✓ جفت کردن روش‌های انقباضی جدید هم‌چون LASSO بیزی، LASSO گروهی، LAD-SCAD

با مدل غیر خطی شبکه عصبی مصنوعی و بررسی کارایی آن‌ها در انتخاب متغیر برای مطالعات QSAR

✓ استفاده از روش‌های انقباضی مورد استفاده در رساله حاضر با روش یادگیری عمیق به‌عنوان روش

مدل‌سازی پر کاربرد و کارآمد

✓ به‌کارگیری روش‌های پیش‌پردازش هم‌چون روش غربالگری مستقل مطمئن (SIS) و الگوریتم تکرار

شونده آن (ISIS) قبل از انتخاب متغیر با روش‌های انقباضی

✓ استفاده از روش‌های انقباضی برای انتخاب متغیر در مطالعات QSAR مربوط به بیماری‌های آلزایمر،

سرطان‌های پروستات و سینه و بررسی میزان کارایی این روش‌ها در مطالعات QSAR این دسته ترکیبات

---

<sup>1</sup>Sure Independence Screening

<sup>2</sup>Iterative Sure Independence Screening

- [1] Young D.C. (2009), "Computational drug design: a guide for computational and medicinal chemists", John Wiley & Sons,
- [2] Wold S. (1995) "Chemometrics; what do we mean with it, and what do we want from it?", Chemometrics and Intelligent Laboratory Systems. **30**, 1, pp 109-115
- [3] Hansch C., and Fujita T. (1964) " $\rho$ - $\sigma$ - $\pi$  Analysis. A method for the correlation of biological activity and chemical structure", Journal of the American Chemical Society. **86**, 8, pp 1616-1626
- [4] Cherkasov A., Muratov E.N., Fourches D., Varnek A., et al. (2014) "QSAR modeling: where have you been? Where are you going to?", Journal of medicinal chemistry. **57**, 12, pp 4977-5010
- [5] Verma J., Khedkar V.M., and Coutinho E.C. (2010) "3D-QSAR in drug design-a review", Current topics in medicinal chemistry. **10**, 1, pp 95-115
- [6] Golbraikh A., Wang X.S., Zhu H., and Tropsha A. (2012) "Predictive QSAR modeling: methods and applications in drug discovery and chemical risk assessment", Handbook of computational chemistry. pp 1309-1342
- [7] Ma J., Sheridan R.P., Liaw A., Dahl G.E., et al. (2015) "Deep neural nets as a method for quantitative structure-activity relationships", Journal of chemical information and modeling. **55**, 2, pp 263-274
- [8] Andrada M.F., Vega-Hissi E.G., Estrada M.R., and Garro Martinez J.C. (2017) "Impact assessment of the rational selection of training and test sets on the predictive ability of QSAR models", SAR and QSAR in Environmental Research. **28**, 12, pp 1011-1023
- [9] Todeschini R., and Consonni V. (2008), "Handbook of molecular descriptors", John Wiley & Sons,
- [10] Tropsha A., Gramatica P., and Gombar V.K. (2003) "The importance of being earnest: validation is the absolute essential for successful application and interpretation of QSPR models", QSAR & Combinatorial Science. **22**, 1, pp 69-77
- [11] Katritzky A.R., Dobchev D.A., Slavov S., and Karelson M. (2008) "Legitimate utilization of large descriptor pools for QSPR/QSAR models", Journal of chemical information and modeling. **48**, 11, pp 2207-2213
- [12] Katritzky A.R., Petrukhin R., Tatham D., Basak S., et al. (2001) "Interpretation of quantitative structure-property and- activity relationships", Journal of chemical information and computer sciences. **41**, 3, pp 679-685
- [13] Jiang J., Duan W., Wei Q., Zhao X., et al. (2020) "Development of quantitative structure-property relationship (QSPR) models for predicting the thermal hazard of ionic liquids: A review of methods and models", Journal of Molecular Liquids. **301**, pp 112471
- [14] Snodin D.J. (2002) "An EU perspective on the use of in vitro methods in regulatory pharmaceutical toxicology", Toxicology letters. **127**, 1-3, pp 161-168
- [15] Kraljevic S., Stambrook P.J., and Pavelic K. (2004) "Accelerating drug discovery: Although the evolution of 'omics' methodologies is still in its infancy, both the pharmaceutical industry and patients could benefit from their implementation in the drug development process", EMBO reports. **5**, 9, pp 837-842
- [16] Adams C.P., and Brantner V.V. (2006) "Estimating the cost of new drug development:

- is it really \$802 million?", *Health affairs*. **25**, 2, pp 420-428
- [17] Kola I., and Landis J. (2004) "Can the pharmaceutical industry reduce attrition rates?", *Nature reviews Drug discovery*. **3**, 8, pp 711-716
- [18] Lionberger R.A. (2008) "FDA critical path initiatives: opportunities for generic drug development", *The AAPS journal*. **10**, 1, pp 103-109
- [19] Andricopulo A.D., Guido R.V., and Oliva G. (2008) "Virtual screening and its integration with modern drug design technologies", *Current medicinal chemistry*. **15**, 1, pp 37-46
- [20] Schwaighofer A., Schroeter T., Mika S., and Blanchard G. (2009) "How wrong can we get? A review of machine learning approaches and error bars", *Combinatorial chemistry & high throughput screening*. **12**, 5, pp 453-468
- [21] Valerio Jr L.G. (2009) "In silico toxicology for the pharmaceutical sciences", *Toxicology and applied pharmacology*. **241**, 3, pp 356-370
- [22] Yap C., Xue Y., and Chen Y. (2006) "Application of support vector machines to in silico prediction of cytochrome P450 enzyme substrates and inhibitors", *Current topics in medicinal chemistry*. **6**, 15, pp 1593-1607
- [23] Meeting of the Chemicals Committee OECD principles, For the validation of (quantitative) structure–activity relationship models, **2019**
- [24] Hdoufane I., Ounaissi D., Dermoune A., and Cherqaoui D. (2021) "Development of QSAR Models Using Singular Value Decomposition Method: A Case Study for Predicting Anti-HIV-1 and Anti-HCV Biological Activities", *Biointerface Research in Applied Chemistry*. **12**, 3, pp 3090-3105
- [25] Puzyn T., Mostrag-Szlichtyng A., Gajewicz A., Skrzyński M., et al. (2011) "Investigating the influence of data splitting on the predictive ability of QSAR/QSPR models", *Structural Chemistry*. **22**, 4, pp 795-804
- [26] HyperChem 8 (by Hypercube, Inc.), Coleman W.F., and Arumainayagam C.R., **2012**
- [27] R: A language and environment for statistical computing, Team R.C., Vienna, Austria, **2020**
- [28] Package ‘caret’, Kuhn M., Wing J., Weston S., Williams A., et al., *The R Journal*, **2020**
- [29] Tibshirani R. (1996) "Regression shrinkage and selection via the lasso", *Journal of the Royal Statistical Society: Series B (Methodological)*. **58**, 1, pp 267-288
- [30] Zou H. (2006) "The adaptive lasso and its oracle properties", *Journal of the American statistical association*. **101**, 476, pp 1418-1429
- [31] Fan J., and Li R. (2001) "Variable selection via nonconcave penalized likelihood and its oracle properties", *Journal of the American statistical Association*. **96**, 456, pp 1348-1360
- [32] Wang H., Li G., and Jiang G. (2007) "Robust regression shrinkage and consistent variable selection through the LAD-Lasso", *Journal of Business & Economic Statistics*. **25**, 3, pp 347-355
- [33] Package ‘ncvreg’, Breheny P., and Breheny M.P., **2021**
- [34] Package ‘parcor’, Kraemer N., Schaefer J., and Kraemer M.N., **2014**
- [35] Package ‘quantreg’, Koenker R., Portnoy S., Ng P.T., Zeileis A., et al., *Cran R-project.org*, **2018**
- [36] Jammalamadaka S.R., *Introduction to linear regression analysis*, in, Taylor & Francis, 2003.
- [37] Saleh A.M.E., Arashi M., and Kibria B.G. (2019), "Theory of ridge regression estimation with applications", *John Wiley & Sons*,

- [38] Sadeghi F., Afkhami A., Madrakian T., and Ghavami R. (2021) "A new approach for simultaneous calculation of pIC<sub>50</sub> and logP through QSAR/QSPR modeling on anthracycline derivatives: a comparable study", *Journal of the Iranian Chemical Society*. pp 1-16
- [39] Chamjangali M.A., Beglari M., and 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", *Journal of Molecular Graphics and Modelling*. **26**, 1, pp 360-367
- [40] Zupan J., and Gasteiger J. (1993), "Neural networks for chemists; an introduction", VCH publishers,
- [41] Bring J. (1994) "How to standardize regression coefficients", *The American Statistician*. **48**, 3, pp 209-213
- [42] Tropsha A. (2010) "Best practices for QSAR model development, validation, and exploitation", *Molecular informatics*. **29**, 6- 7, pp 476-488
- [43] Golbraikh A., and Tropsha A. (2002) "Beware of q<sup>2</sup>!", *Journal of molecular graphics and modelling*. **20**, 4, pp 269-276
- [44] Roy P.P., and Roy K. (2008) "On some aspects of variable selection for partial least squares regression models", *QSAR & Combinatorial Science*. **27**, 3, pp 302-313
- [45] Consonni V., Ballabio D., and Todeschini R. (2009) "Comments on the definition of the Q<sup>2</sup> parameter for QSAR validation", *Journal of chemical information and modeling*. **49**, 7, pp 1669-1678
- [46] Consonni V., Ballabio D., and Todeschini R. (2010) "Evaluation of model predictive ability by external validation techniques", *Journal of chemometrics*. **24**, 3- 4, pp 194-201
- [47] Alonso H., Bliznyuk A.A., and Gready J.E. (2006) "Combining docking and molecular dynamic simulations in drug design", *Medicinal research reviews*. **26**, 5, pp 531-568
- [48] Wermuth C., Ganellin C., Lindberg P., and Mitscher L. (1998) "Glossary of terms used in medicinal chemistry (IUPAC Recommendations 1998)", *Pure and Applied Chemistry*. **70**, 5, pp 1129-1143
- [49] Liao C., Peach M.L., Yao R., and Nicklaus M.C., *Molecular docking and structure-based virtual screening*, in, *Future Medicine*, **2013**
- [50] Feinstein W.P., and Brylinski M. (2015) "Calculating an optimal box size for ligand docking and virtual screening against experimental and predicted binding pockets", *Journal of cheminformatics*. **7**, 1, pp 1-10
- [51] Goodsell D.S., Morris G.M., and Olson A.J. (1996) "Automated docking of flexible ligands: applications of AutoDock", *Journal of molecular recognition*. **9**, 1, pp 1-5
- [52] Lang P.T., Aynechi T., Moustakas D., Shoichet B., et al. (2007) "Molecular docking and structure-based design", *Drug Discovery Research: New Frontiers in the Post-Genomic Era*. pp 3-23
- [53] Mehellou Y., and De Clercq E. (2010) "Twenty-six years of anti-HIV drug discovery: where do we stand and where do we go?", *Journal of medicinal chemistry*. **53**, 2, pp 521-538
- [54] Namasivayam V., Vanangamudi M., Kramer V.G., Kurup S., et al. (2018) "The journey of HIV-1 non-nucleoside reverse transcriptase inhibitors (NNRTIs) from lab to clinic", *Journal of medicinal chemistry*. **62**, 10, pp 4851-4883
- [55] Zhan P., Pannecouque C., De Clercq E., and Liu X. (2016) "Anti-HIV drug discovery and development: current innovations and future trends: miniperspective", *Journal of*

medicinal chemistry. **59**, 7, pp 2849-2878

[56] Collaboration A.T.C. (2008) "Life expectancy of individuals on combination antiretroviral therapy in high-income countries: a collaborative analysis of 14 cohort studies", *The Lancet*. **372**, 9635, pp 293-299

[57] Brechtel J.R., Breitbart W., Galiotta M., Krivo S., et al. (2001) "The use of highly active antiretroviral therapy (HAART) in patients with advanced HIV infection: impact on medical, palliative care, and quality of life outcomes", *Journal of pain and symptom management*. **21**, 1, pp 41-51

[58] Meanwell N.A., Krystal M.R., Nowicka-Sans B., Langley D.R., et al., *Inhibitors of HIV-1 attachment: the discovery and development of temsavir and its prodrug fostemsavir*, in, ACS Publications, 2018.

[59] De Clercq E. (1998) "The role of non-nucleoside reverse transcriptase inhibitors (NNRTIs) in the therapy of HIV-1 infection", *Antiviral research*. **38**, 3, pp 153-179

[60] De Clercq E. (2004) "Non- nucleoside reverse transcriptase inhibitors (NNRTIs): past, present, and future", *Chemistry & biodiversity*. **1**, 1, pp 44-64

[61] De Clercq E., and Li G. (2016) "Approved antiviral drugs over the past 50 years", *Clinical microbiology reviews*. **29**, 3, pp 695-747

[62] Wu F., Zhao S., Yu B., Chen Y.-M., et al. (2020) "A new coronavirus associated with human respiratory disease in China", *Nature*. **579**, 7798, pp 265-269

[63] Zhou P., Yang X.-L., Wang X.-G., Hu B., et al. (2020) "A pneumonia outbreak associated with a new coronavirus of probable bat origin", *nature*. **579**, 7798, pp 270-273

[64] Fahmi I. (2019) "World Health Organization coronavirus disease 2019 (Covid-19) situation report", *DroneEmprit*. pp 1-9

[65] Azhar E.I., Hui D.S., Memish Z.A., Drosten C., et al. (2019) "The middle east respiratory syndrome (MERS)", *Infectious Disease Clinics*. **33**, 4, pp 891-905

[66] Hui D.S., and Zumla A. (2019) "Severe acute respiratory syndrome: historical, epidemiologic, and clinical features", *Infectious Disease Clinics*. **33**, 4, pp 869-889

[67] Wang C., Horby P.W., Hayden F.G., and Gao G.F. (2020) "A novel coronavirus outbreak of global health concern", *The lancet*. **395**, 10223, pp 470-473

[68] Ciotti M., Ciccozzi M., Terrinoni A., Jiang W.-C., et al. (2020) "The COVID-19 pandemic", *Critical reviews in clinical laboratory sciences*. **57**, 6, pp 365-388

[69] Zhang L., Lin D., Sun X., Curth U., et al. (2020) "Crystal structure of SARS-CoV-2 main protease provides a basis for design of improved  $\alpha$ -ketoamide inhibitors", *Science*. **368**, 6489, pp 409-412

[70] Lee C.-C., Kuo C.-J., Hsu M.-F., Liang P.-H., et al. (2007) "Structural basis of mercury- and zinc-conjugated complexes as SARS-CoV 3C-like protease inhibitors", *FEBS letters*. **581**, 28, pp 5454-5458

[71] Ramajayam R., Tan K.-P., and Liang P.-H. (2011) "Recent development of 3C and 3CL protease inhibitors for anti-coronavirus and anti-picornavirus drug discovery", *Biochemical Society Transactions*. **39**, 5, pp 1371-1375

[72] Barta I., Smerak P., Polivkova Z., Sestakova H., et al. (2006) "Current trends and perspectives in nutrition and cancer prevention", *Neoplasma*. **53**, 1, pp 19-25

[73] Siegel R.L., Miller K.D., and Jemal A. (2020) "Cancer statistics, 2020", *CA: a cancer journal for clinicians*. **70**, 1, pp 7-30

[74] Bade B.C., and Cruz C.S.D. (2020) "Lung cancer 2020: epidemiology, etiology, and prevention", *Clinics in chest medicine*. **41**, 1, pp 1-24

- [75] Garcia-Echeverria C., and Sellers W. (2008) "Drug discovery approaches targeting the PI3K/Akt pathway in cancer", *Oncogene*. **27**, 41, pp 5511-5526
- [76] Ran T., Lu T., Yuan H., Liu H., et al. (2012) "A selectivity study on mTOR/PI3K $\alpha$  inhibitors by homology modeling and 3D-QSAR", *Journal of molecular modeling*. **18**, 1, pp 171-186
- [77] Walker E.H., Pacold M.E., Perisic O., Stephens L., et al. (2000) "Structural determinants of phosphoinositide 3-kinase inhibition by wortmannin, LY294002, quercetin, myricetin, and staurosporine", *Molecular cell*. **6**, 4, pp 909-919
- [78] Arcaro A., Volinia S., Zvelebil M.J., Stein R., et al. (1998) "Human phosphoinositide 3-kinase C2 $\beta$ , the role of calcium and the C2 domain in enzyme activity", *Journal of Biological Chemistry*. **273**, 49, pp 33082-33090
- [79] Sadhu C., Masinovsky B., Dick K., Sowell C.G., et al. (2003) "Essential role of phosphoinositide 3-kinase  $\delta$  in neutrophil directional movement", *The Journal of Immunology*. **170**, 5, pp 2647-2654
- [80] Yin Y., Wu X., Han H.-W., Sha S., et al. (2014) "Discovery and synthesis of a novel series of potent, selective inhibitors of the PI3K $\alpha$ : 2-alkyl-chromeno [4, 3-c] pyrazol-4 (2 H)-one derivatives", *Organic & biomolecular chemistry*. **12**, 45, pp 9157-9165
- [81] Lu L., Sha S., Wang K., Zhang Y.-H., et al. (2016) "Discovery of chromeno [4, 3-c] pyrazol-4 (2H)-one containing carbonyl or oxime derivatives as potential, selective inhibitors PI3K $\alpha$ ", *Chemical and Pharmaceutical Bulletin*. pp c16-00388
- [82] Sin D.W.-m., Wong Y.-c., Sham W.-c., and Wang D. (2001) "Development of an analytical technique and stability evaluation of 143 C3–C12 volatile organic compounds in Summa® canisters by gas chromatography–mass spectrometry", *Analyst*. **126**, 3, pp 310-321
- [83] Luan F., Xue C., Zhang R., Zhao C., et al. (2005) "Prediction of retention time of a variety of volatile organic compounds based on the heuristic method and support vector machine", *Analytica Chimica Acta*. **537**, 1-2, pp 101-110
- [84] Jalali-Heravi M., and Kyani A. (2004) "Use of computer-assisted methods for the modeling of the retention time of a variety of volatile organic compounds: a PCA-MLR-ANN approach", *Journal of chemical information and computer sciences*. **44**, 4, pp 1328-1335
- [85] Sepehri B., Ghavami R., Farahbakhsh S., and Ahmadi R. (2021) "Machine learning-based quantitative structure–retention relationship models for predicting the retention indices of volatile organic pollutants", *International Journal of Environmental Science and Technology*. pp 1-10
- [86] Zhu T., and Tao C. (2022) "Prediction models with multiple machine learning algorithms for POPs: The calculation of PDMS-air partition coefficient from molecular descriptor", *Journal of Hazardous Materials*. **423**, pp 127037
- [87] Paskaleva V., and Kochev N. (2019) "GROUP CONTRIBUTION MODELING OF RETENTION INDICES OF VOLATILE ORGANIC COMPOUNDS IN PEPPERS", *Scientific Work of the*. pp 243
- [88] He M., Yan P., Yang Z., Ye Y., et al. (2018) "Multi-analytical strategy for unassigned peaks using physical/mathematical separation, fragmental rules and retention index prediction: An example of sesquiterpene metabolites characterization in *Cyperus rotundus*", *Journal of pharmaceutical and biomedical analysis*. **154**, pp 476-485
- [89] Kaliszan R. (2007) "QSRR: quantitative structure-(chromatographic) retention



relationships", *Chemical reviews*. **107**, 7, pp 3212-3246

[90] Kalisz R., and Bączek T., *QSAR in chromatography: quantitative structure–retention relationships (QSRRs)*, in: *Recent advances in QSAR studies*, Springer, **2010**

[91] Marrero-Ponce Y., Barigye S.J., Jorge-Rodríguez M.E., and Tran-Thi-Thu T. (2018) "QSRR prediction of gas chromatography retention indices of essential oil components", *Chemical Papers*. **72**, 1, pp 57-69

[92] Fragkaki A., Tsantili-Kakoulidou A., Angelis Y., Koupparis M., et al. (2009) "Gas chromatographic quantitative structure–retention relationships of trimethylsilylated anabolic androgenic steroids by multiple linear regression and partial least squares", *Journal of Chromatography A*. **1216**, 47, pp 8404-8420

[93] Jalali-Heravi M., and Fatemi M.H. (2001) "Artificial neural network modeling of Kovats retention indices for noncyclic and monocyclic terpenes", *Journal of Chromatography A*. **915**, 1-2, pp 177-183

[94] Golmohammadi H., and Fatemi M.H. (2005) "Artificial neural network prediction of retention factors of some benzene derivatives and heterocyclic compounds in micellar electrokinetic chromatography", *Electrophoresis*. **26**, 18, pp 3438-3444

[95] Algamal Z., and Lee M. (2017) "A new adaptive L1-norm for optimal descriptor selection of high-dimensional QSAR classification model for anti-hepatitis C virus activity of thiourea derivatives", *SAR and QSAR in Environmental Research*. **28**, 1, pp 75-90

[96] Algamal Z., Qasim M., and Ali H. (2017) "A QSAR classification model for neuraminidase inhibitors of influenza A viruses (H1N1) based on weighted penalized support vector machine", *SAR and QSAR in Environmental Research*. **28**, 5, pp 415-426

[97] Algamal Z.Y., Alhamzawi R., and Ali H.T.M. (2018) "Gene selection for microarray gene expression classification using Bayesian Lasso quantile regression", *Computers in biology and medicine*. **97**, pp 145-152

[98] Algamal Z.Y., and Lee M.H. (2017) "A novel molecular descriptor selection method in QSAR classification model based on weighted penalized logistic regression", *Journal of Chemometrics*. **31**, 10, pp e2915

[99] Al-Thanoon N.A., Qasim O.S., and Algamal Z.Y. (2019) "A new hybrid firefly algorithm and particle swarm optimization for tuning parameter estimation in penalized support vector machine with application in chemometrics", *Chemometrics and Intelligent Laboratory Systems*. **184**, pp 142-152

[100] Qasim M., Algamal Z., and Ali H.M. (2018) "A binary QSAR model for classifying neuraminidase inhibitors of influenza A viruses (H1N1) using the combined minimum redundancy maximum relevancy criterion with the sparse support vector machine", *SAR and QSAR in Environmental Research*. **29**, 7, pp 517-527

[101] Qasim O.S., Al-Thanoon N.A., and Algamal Z.Y. (2020) "Feature selection based on chaotic binary black hole algorithm for data classification", *Chemometrics and Intelligent Laboratory Systems*. **204**, pp 104104

[102] Wu S., Jiang H., Shen H., and Yang Z. (2018) "Gene selection in cancer classification using sparse logistic regression with L1/2 regularization", *Applied Sciences*. **8**, 9, pp 1569

[103] Algamal Z.Y., Lee M.H., Al-Fakih A.M., and Aziz M. (2017) "High- dimensional QSAR classification model for anti- hepatitis C virus activity of thiourea derivatives based on the sparse logistic regression model with a bridge penalty", *Journal of Chemometrics*. **31**, 6, pp e2889

[104] Alharthi A., Lee M., Algamal Z., and Al-Fakih A. (2020) "Quantitative structure-

activity relationship model for classifying the diverse series of antifungal agents using ratio weighted penalized logistic regression", SAR and QSAR in Environmental Research. **31**, 8, pp 571-583

[105] Alharthi A.M., Lee M.H., and Algamal Z.Y. (2021) "Gene selection and classification of microarray gene expression data based on a new adaptive L1-norm elastic net penalty", Informatics in Medicine Unlocked. pp 100622

[106] Algamal Z. (2017) "An efficient gene selection method for high-dimensional microarray data based on sparse logistic regression", Electronic Journal of Applied Statistical Analysis. **10**, 1, pp 242-256

[107] Peng X.-L., Yin H., Li R., and Fang K.-T. (2006) "The application of Kriging and empirical Kriging based on the variables selected by SCAD", Analytica chimica acta. **578**, 2, pp 178-185

[108] Algamal Z.Y., Lee M.H., Al-Fakih A.M., and Aziz M. (2015) "High-dimensional QSAR prediction of anticancer potency of imidazo [4, 5- b] pyridine derivatives using adjusted adaptive LASSO", Journal of Chemometrics. **29**, 10, pp 547-556

[109] Algamal Z., Lee M., Al-Fakih A., and Aziz M. (2016) "High-dimensional QSAR modelling using penalized linear regression model with L 1/2-norm", SAR and QSAR in Environmental Research. **27**, 9, pp 703-719

[110] Al-Fakih A., Algamal Z., Lee M., and Aziz M. (2018) "A penalized quantitative structure–property relationship study on melting point of energetic carbocyclic nitroaromatic compounds using adaptive bridge penalty", SAR and QSAR in Environmental Research. **29**, 5, pp 339-353

[111] Majumdar S., Basak S.C., Lungu C., Diudea M., et al. (2018) "Mathematical structural descriptors and mutagenicity assessment: a study with congeneric and diverse datasets", SAR and QSAR in Environmental Research. **29**, 8, pp 579-590

[112] Xia L.-Y., Wang Y.-W., Meng D.-Y., Yao X.-J., et al. (2018) "Descriptor selection via log-sum regularization for the biological activities of chemical structure", International journal of molecular sciences. **19**, 1, pp 30

[113] Al-Dabbagh Z.T., and Algamal Z.Y. (2019) "Least absolute deviation estimator-bridge variable selection and estimation for quantitative structure–activity relationship model", Journal of Chemometrics. **33**, 7, pp e3139

[114] Mozafari Z., Arab Chamjangali M., Arashi M., and Goudarzi N. (2021) "Performance of smoothly clipped absolute deviation as a variable selection method in the artificial neural network- based QSAR studies", Journal of Chemometrics. **35**, 5, pp e3338

[115] Mozafari Z., Chamjangali M.A., Arashi M., and Goudarzi N. (2021) "Suggestion of active 3-chymotrypsin like protease (3CLPro) inhibitors as potential anti-SARS-CoV-2 agents using predictive QSAR model based on the combination of ALASSO with an ANN model", SAR and QSAR in Environmental Research. **32**, 11, pp 863-888

[116] HyperChem(TM) Professional 8.0.5, Hypercube I., 1115 NW 4th Street, Gainesville, Florida 32601, USA., 2008

[117] Dragon software: An easy approach to molecular descriptor calculations, Mauri A., Consonni V., Pavan M., and Todeschini R., University of Milano-Bicocca, 2006

[118] IBM SPSS statistics 25 for Windows student, George D., and Mallery P., 2017

[119] Matlab Version R2017a, The MathWorks Inc., Massachusetts, 2017

[120] Origin Lab Version 2021 , OriginLab Corporation, Northampton, MA, USA., 2021

- [121] Morris G.M., Huey R., Lindstrom W., Sanner M.F., et al. (2009) "AutoDock4 and AutoDockTools4: Automated docking with selective receptor flexibility", *Journal of computational chemistry*. **30**, 16, pp 2785-2791
- [122] ViewerLite v 5.0, Accelrys Inc, San Diego, CA, 2001
- [123] BIOVIA discovery studio, Biovia D.S., San Diego, CA, USA, 2021
- [124] Moss J.A. (2013) "HIV/AIDS Review", *Radiologic technology*. **84**, 3, pp 247-267
- [125] de Béthune M.-P. (2010) "Non-nucleoside reverse transcriptase inhibitors (NNRTIs), their discovery, development, and use in the treatment of HIV-1 infection: a review of the last 20 years (1989–2009)", *Antiviral research*. **85**, 1, pp 75-90
- [126] Reust C.E. (2011) "Common adverse effects of antiretroviral therapy for HIV disease", *American family physician*. **83**, 12, pp 1443-1451
- [127] Chen X., Liu X., Meng Q., Wang D., et al. (2013) "Novel piperidinylamino-diarylpyrimidine derivatives with dual structural conformations as potent HIV-1 non-nucleoside reverse transcriptase inhibitors", *Bioorganic & medicinal chemistry letters*. **23**, 24, pp 6593-6597
- [128] De La Rosa M., Kim H.W., Gunic E., Jenket C., et al. (2006) "Tri-substituted triazoles as potent non-nucleoside inhibitors of the HIV-1 reverse transcriptase", *Bioorganic & medicinal chemistry letters*. **16**, 17, pp 4444-4449
- [129] Gagnon A., Landry S., Coulombe R., Jakalian A., et al. (2009) "Investigation on the role of the tetrazole in the binding of thiotetrazolylacetanilides with HIV-1 wild type and K103N/Y181C double mutant reverse transcriptases", *Bioorganic & medicinal chemistry letters*. **19**, 4, pp 1199-1205
- [130] Moyle G., Boffito M., Stoehr A., Rieger A., et al. (2010) "Phase 2a randomized controlled trial of short-term activity, safety, and pharmacokinetics of a novel nonnucleoside reverse transcriptase inhibitor, RDEA806, in HIV-1-positive, antiretroviral-naive subjects", *Antimicrobial agents and chemotherapy*. **54**, 8, pp 3170-3178
- [131] O'Meara J.A., Jakalian A., LaPlante S., Bonneau P.R., et al. (2007) "Scaffold hopping in the rational design of novel HIV-1 non-nucleoside reverse transcriptase inhibitors", *Bioorganic & medicinal chemistry letters*. **17**, 12, pp 3362-3366
- [132] Tzoupis H., Leonis G., Durdagi S., Mouchlis V., et al. (2011) "Binding of novel fullerene inhibitors to HIV-1 protease: insight through molecular dynamics and molecular mechanics Poisson–Boltzmann surface area calculations", *Journal of computer-aided molecular design*. **25**, 10, pp 959-976
- [133] Wang Z., Wu B., Kuhlen K.L., Bursulaya B., et al. (2006) "Synthesis and biological evaluations of sulfanyltriazoles as novel HIV-1 non-nucleoside reverse transcriptase inhibitors", *Bioorganic & medicinal chemistry letters*. **16**, 16, pp 4174-4177
- [134] Li X., Lu X., Chen W., Liu H., et al. (2014) "Arylazolyl (azinyl) thioacetanilides. Part 16: Structure-based bioisosterism design, synthesis and biological evaluation of novel pyrimidinylthioacetanilides as potent HIV-1 inhibitors", *Bioorganic & medicinal chemistry*. **22**, 19, pp 5290-5297
- [135] Zhan P., Chen W., Li Z., Li X., et al. (2012) "Discovery of novel 2-(3-(2-chlorophenyl)pyrazin-2-ylthio)-N-arylacetamides as potent HIV-1 inhibitors using a structure-based bioisosterism approach", *Bioorganic & medicinal chemistry*. **20**, 23, pp 6795-6802
- [136] Zhan P., Li X., Li Z., Chen X., et al. (2012) "Structure-based bioisosterism design, synthesis and biological evaluation of novel 1, 2, 4-triazin-6-ylthioacetamides as potent HIV-1 NNRTIs", *Bioorganic & medicinal chemistry letters*. **22**, 23, pp 7155-7162

- [137] Armand-Ugón M., Moncunill G., Gonzalez E., Mena M., et al. (2010) "Different selection patterns of resistance and cross-resistance to HIV-1 agents targeting CCR5", *Journal of antimicrobial chemotherapy*. **65**, 3, pp 417-424
- [138] Beiknejad D., Chaichi M.J., and Fatemi M.H. (2021) "Prediction of ozonolytic decolorization half-lives of azo dyes in a continuous-flow system using QSPR", *Dyes and Pigments*. **185**, pp 108915
- [139] Montgomery D.C., Peck E.A., and Vining G.G. (2021), "Introduction to linear regression analysis", John Wiley & Sons,
- [140] Coulibaly P., Anctil F., and Bobée B. (1999) "Prévision hydrologique par réseaux de neurones artificiels: état de l'art", *Canadian Journal of civil engineering*. **26**, 3, pp 293-304
- [141] Othman F., and Naseri M. (2011) "Reservoir inflow forecasting using artificial neural network", *International journal of physical sciences*. **6**, 3, pp 434-440
- [142] Sahigara F., Mansouri K., Ballabio D., Mauri A., et al. (2012) "Comparison of different approaches to define the applicability domain of QSAR models", *Molecules*. **17**, 5, pp 4791-4810
- [143] Tsang K.W., Ho P.L., Ooi G.C., Yee W.K., et al. (2003) "A cluster of cases of severe acute respiratory syndrome in Hong Kong", *New England Journal of Medicine*. **348**, 20, pp 1977-1985
- [144] Lu H., Stratton C.W., and Tang Y.W. (2020) "Outbreak of pneumonia of unknown etiology in Wuhan, China: The mystery and the miracle", *Journal of medical virology*. **92**, 4, pp 401
- [145] Zhu N., Zhang D., Wang W., Li X., et al. (2020) "A novel coronavirus from patients with pneumonia in China, 2019", *New England journal of medicine*. pp 1-10
- [146] Sohrabi C., Alsafi Z., O'Neill N., Khan M., et al. (2020) "World Health Organization declares global emergency: A review of the 2019 novel coronavirus (COVID-19)", *International journal of surgery*. **76**, pp 71-76
- [147] Chen Y.W., Yiu C.-P.B., and Wong K.-Y. (2020) "Prediction of the SARS-CoV-2 (2019-nCoV) 3C-like protease (3CL pro) structure: virtual screening reveals velpatasvir, ledipasvir, and other drug repurposing candidates", *F1000Research*. **9**, pp 10-16
- [148] Zhavoronkov A., Zagribelnyy B., Zhebrak A., Aladinskiy V., et al. (2020) "Potential non-covalent SARS-CoV-2 3C-like protease inhibitors designed using generative deep learning approaches and reviewed by human medicinal chemist in virtual reality", pp 13-19
- [149] Chen L.-R., Wang Y.-C., Lin Y.W., Chou S.-Y., et al. (2005) "Synthesis and evaluation of isatin derivatives as effective SARS coronavirus 3CL protease inhibitors", *Bioorganic & Medicinal Chemistry Letters*. **15**, 12, pp 3058-3062
- [150] Liu W., Zhu H.-M., Niu G.-J., Shi E.-Z., et al. (2014) "Synthesis, modification and docking studies of 5-sulfonyl isatin derivatives as SARS-CoV 3C-like protease inhibitors", *Bioorganic & medicinal chemistry*. **22**, 1, pp 292-302
- [151] Chen C.-N., Lin C.P., Huang K.-K., Chen W.-C., et al. (2005) "Inhibition of SARS-CoV 3C-like protease activity by theaflavin-3, 3'-digallate (TF3)", *Evidence-Based Complementary and Alternative Medicine*. **2**, 2, pp 209-215
- [152] Chen L., Gui C., Luo X., Yang Q., et al. (2005) "Cinanserin is an inhibitor of the 3C-like proteinase of severe acute respiratory syndrome coronavirus and strongly reduces virus replication in vitro", *Journal of virology*. **79**, 11, pp 7095-7103
- [153] Niu C., Yin J., Zhang J., Vederas J.C., et al. (2008) "Molecular docking identifies the

binding of 3-chloropyridine moieties specifically to the S1 pocket of SARS-CoV Mpro", *Bioorganic & medicinal chemistry*. **16**, 1, pp 293-302

[154] Lu I.-L., Mahindroo N., Liang P.-H., Peng Y.-H., et al. (2006) "Structure-based drug design and structural biology study of novel nonpeptide inhibitors of severe acute respiratory syndrome coronavirus main protease", *Journal of medicinal chemistry*. **49**, 17, pp 5154-5161

[155] Tsai K.-C., Chen S.-Y., Liang P.-H., Lu I.-L., et al. (2006) "Discovery of a novel family of SARS-CoV protease inhibitors by virtual screening and 3D-QSAR studies", *Journal of medicinal chemistry*. **49**, 12, pp 3485-3495

[156] Park J.-Y., Kim J.H., Kim Y.M., Jeong H.J., et al. (2012) "Tanshinones as selective and slow-binding inhibitors for SARS-CoV cysteine proteases", *Bioorganic & medicinal chemistry*. **20**, 19, pp 5928-5935

[157] Dooley A.J., Shindo N., Taggart B., Park J.-G., et al. (2006) "From genome to drug lead: identification of a small-molecule inhibitor of the SARS virus", *Bioorganic & medicinal chemistry letters*. **16**, 4, pp 830-833

[158] Kao R.Y., Tsui W.H., Lee T.S., Tanner J.A., et al. (2004) "Identification of novel small-molecule inhibitors of severe acute respiratory syndrome-associated coronavirus by chemical genetics", *Chemistry & biology*. **11**, 9, pp 1293-1299

[159] Shehroz M., Zaheer T., and Hussain T. (2020) "Computer-aided drug design against spike glycoprotein of SARS-CoV-2 to aid COVID-19 treatment", *Heliyon*. **6**, 10, pp e05278

[160] Bacha U., Barrila J., Velazquez-Campoy A., Leavitt S.A., et al. (2004) "Identification of novel inhibitors of the SARS coronavirus main protease 3CLpro", *Biochemistry*. **43**, 17, pp 4906-4912

[161] Fujita T., and Winkler D.A. (2016) "Understanding the roles of the "two QSARs"", *Journal of chemical information and modeling*. **56**, 2, pp 269-274

[162] Gramatica P. (2020) "Principles of QSAR modeling: comments and suggestions from personal experience", *International Journal of Quantitative Structure-Property Relationships (IJQSPR)*. **5**, 3, pp 61-97

[163] The caret package. R Foundation for Statistical Computing, Vienna, Austria, Kuhn M., URL <https://cran.r-project.org/package=caret>, 2012

[164] Ebrahimi M., Khayamian T., and Gharaghani S. (2012) "Interactions between Activin-Like Kinase 5 (ALK5) receptor and its inhibitors and the construction of a Docking Descriptor-Based QSAR model", *Journal of the Brazilian Chemical Society*. **23**, 11, pp 2043-2092

[165] Eklund M., Norinder U., Boyer S., and Carlsson L. (2012) "Benchmarking variable selection in QSAR", *Molecular informatics*. **31**, 2, pp 173-179

[166] Ghasemi G., Arshadi S., Rashtehroodi A.N., Nirouei M., et al. (2013) "QSAR investigation on quinolizidinyl derivatives in Alzheimer's disease", *Journal of Computational Medicine*. **2013**, pp 13-18

[167] Wacker S., and Noskov S.Y. (2018) "Performance of machine learning algorithms for qualitative and quantitative prediction drug blockade of hERG1 channel", *Computational Toxicology*. **6**, pp 55-63

[168] Zhu X.-W., Xin Y.-J., and Ge H.-L. (2015) "Recursive random forests enable better predictive performance and model interpretation than variable selection by LASSO", *Journal of chemical information and modeling*. **55**, 4, pp 736-746

[169] Chen W., Zhan P., Rai D., De Clercq E., et al. (2014) "Discovery of 2-pyridone derivatives as potent HIV-1 NNRTIs using molecular hybridization based on

- crystallographic overlays", *Bioorganic & medicinal chemistry*. **22**, 6, pp 1863-1872
- [170] Chen X., Zhan P., Liu X., Cheng Z., et al. (2012) "Design, synthesis, anti-HIV evaluation and molecular modeling of piperidine-linked amino-triazine derivatives as potent non-nucleoside reverse transcriptase inhibitors", *Bioorganic & medicinal chemistry*. **20**, 12, pp 3856-3864
- [171] Li D., Zhan P., Liu H., Pannecouque C., et al. (2013) "Synthesis and biological evaluation of pyridazine derivatives as novel HIV-1 NNRTIs", *Bioorganic & medicinal chemistry*. **21**, 7, pp 2128-2134
- [172] Wang J., Zhan P., Li Z., Liu H., et al. (2014) "Discovery of nitropyridine derivatives as potent HIV-1 non-nucleoside reverse transcriptase inhibitors via a structure-based core refining approach", *European journal of medicinal chemistry*. **76**, pp 531-538
- [173] Yin Y., Hu J.-Q., Wu X., Sha S., et al. (2019) "Design, synthesis and biological evaluation of novel chromeno [4, 3-c] pyrazol-4 (2H)-one derivatives containing sulfonamido as potential PI3K $\alpha$  inhibitors", *Bioorganic & medicinal chemistry*. **27**, 11, pp 2261-2267
- [174] Borhani T.N., García-Muñoz S., Luciani C.V., Galindo A., et al. (2019) "Hybrid QSPR models for the prediction of the free energy of solvation of organic solute/solvent pairs", *Physical Chemistry Chemical Physics*. **21**, 25, pp 13706-13720
- [175] Li M., Yu H., Wang Y., Li J., et al. (2020) "QSPR models for predicting the adsorption capacity for microplastics of polyethylene, polypropylene and polystyrene", *Scientific reports*. **10**, 1, pp 1-11
- [176] Raffetti E., Donato F., Speziani F., Scarcella C., et al. (2018) "Polychlorinated biphenyls (PCBs) exposure and cardiovascular, endocrine and metabolic diseases: a population-based cohort study in a North Italian highly polluted area", *Environment international*. **120**, pp 215-222
- [177] Sepehri B. (2020) "A review on created QSPR models for predicting ionic liquids properties and their reliability from chemometric point of view", *Journal of Molecular Liquids*. **297**, pp 112013
- [178] Villaverde J.J., Sevilla-Morán B., López-Goti C., Alonso-Prados J.L., et al. (2018) "Considerations of nano-QSAR/QSPR models for nanopesticide risk assessment within the European legislative framework", *Science of the Total Environment*. **634**, pp 1530-1539
- [179] Yan F., Shi Y., Wang Y., Jia Q., et al. (2020) "QSPR models for the properties of ionic liquids at variable temperatures based on norm descriptors", *Chemical Engineering Science*. **217**, pp 115540
- [180] Zhu T., Gu L., Chen M., and Sun F. (2021) "Exploring QSPR models for predicting PUF-air partition coefficients of organic compounds with linear and nonlinear approaches", *Chemosphere*. **266**, pp 128962
- [181] Acimovic M., Pezo L., Jeremic J.S., Cvetkovic M., et al. (2020) "QSRR model for predicting retention indices of geraniol chemotype of *Thymus serpyllum* essential oil", *Journal of Essential Oil Bearing Plants*. **23**, 3, pp 464-473
- [182] Kaliszan R. (2020) "Recent advances in quantitative structure-retention relationships", *Handbook of Analytical Separations*. **8**, pp 587-632
- [183] Matyushin D.D., Sholokhova A.Y., Karnaeva A.E., and Buryak A.K. (2020) "Various aspects of retention index usage for GC-MS library search: A statistical investigation using a diverse data set", *Chemometrics and Intelligent Laboratory Systems*. **202**, pp 104042
- [184] Naylor B.C., Catrow J.L., Maschek J.A., and Cox J.E. (2020) "QSRR automator: a

tool for automating retention time prediction in lipidomics and metabolomics", *Metabolites*. **10**, 6, pp 237

[185] Pavlić B., Teslić N., Kojić P., and Pezo L. (2020) "Prediction of the GC-MS retention time for terpenoids detected in sage (*Salvia officinalis* L.) essential oil using QSRR approach", *Journal of the Serbian Chemical Society*. **85**, 1, pp 9-23

[186] Wu L., Gong P., Wu Y., Liao K., et al. (2013) "An integral strategy toward the rapid identification of analogous nontarget compounds from complex mixtures", *Journal of Chromatography A*. **1303**, pp 39-47

[187] Kalhor P., and Yarivand O. (2016) "Quantitative Structure-Property Relationship Study to Predict the Retention Times of Some Volatile Compounds in Rosé Wines", *Analytical Chemistry Letters*. **6**, 4, pp 371-383

[188] Andries J.P., Goodarzi M., and Vander Heyden Y. (2020) "Improvement of quantitative structure-retention relationship models for chromatographic retention prediction of peptides applying individual local partial least squares models", *Talanta*. **219**, pp 121266

[189] Zheng W., and Jin M. (2020) "Comparing multiple categories of feature selection methods for text classification", *Digital Scholarship in the Humanities*. **35**, 1, pp 208-224

[190] Al-Fakih A., Algamal Z., Lee M., and Aziz M. (2017) "A sparse QSRR model for predicting retention indices of essential oils based on robust screening approach", *SAR and QSAR in Environmental Research*. **28**, 8, pp 691-703

[191] Fleming-Jones M.E., and Smith R.E. (2003) "Volatile organic compounds in foods: a five year study", *Journal of agricultural and food chemistry*. **51**, 27, pp 8120-8127

[192] Vinci R.M., Jacxsens L., De Meulenaer B., Deconink E., et al. (2015) "Occurrence of volatile organic compounds in foods from the Belgian market and dietary exposure assessment", *Food Control*. **52**, pp 1-8

[193] Ghavami R., and Faham S. (2010) "QSRR models for Kováts' retention indices of a variety of volatile organic compounds on polar and apolar GC stationary phases using molecular connectivity indexes", *Chromatographia*. **72**, 9, pp 893-903

[194] Xu J., Zhang W., Adhikari K., and Shi Y.-C. (2017) "Determination of volatile compounds in heat-treated straight-grade flours from normal and waxy wheats", *Journal of Cereal Science*. **75**, pp 77-83

[195] Häppölä P., Havulinna A.S., Tasa T., Mars N.J., et al. (2020) "A data-driven medication score predicts 10-year mortality among aging adults", *Scientific reports*. **10**, 1, pp 1-10

[196] Ren J., Chamberlain P.P., Stamp A., Short S.A., et al. (2008) "Structural basis for the improved drug resistance profile of new generation benzophenone non-nucleoside HIV-1 reverse transcriptase inhibitors", *Journal of medicinal chemistry*. **51**, 16, pp 5000-5008

[197] Hevener K.E., Zhao W., Ball D.M., Babaoglu K., et al. (2009) "Validation of molecular docking programs for virtual screening against dihydropteroate synthase", *Journal of chemical information and modeling*. **49**, 2, pp 444-460

[198] Burley S.K., Berman H.M., Bhikadiya C., Bi C., et al. (2019) "RCSB Protein Data Bank: biological macromolecular structures enabling research and education in fundamental biology, biomedicine, biotechnology and energy", *Nucleic acids research*. **47**, D1, pp D464-D474

[199] Zhou Q., Zhang N., Zhang C., Huang L., et al. (2010) "Molecular mechanism of enantioselective inhibition of acetolactate synthase by imazethapyr enantiomers", *Journal of agricultural and food chemistry*. **58**, 7, pp 4202-4206

- [200] Daina A., Michielin O., and Zoete V. (2017) "SwissADME: a free web tool to evaluate pharmacokinetics, drug-likeness and medicinal chemistry friendliness of small molecules", *Scientific reports*. **7**, 1, pp 1-13
- [201] Alam S., and Khan F. (2014) "QSAR and docking studies on xanthone derivatives for anticancer activity targeting DNA topoisomerase II $\alpha$ ", *Drug design, development and therapy*. **8**, pp 183
- [202] Jin Z., Du X., Xu Y., Deng Y., et al. (2020) "Structure of M pro from SARS-CoV-2 and discovery of its inhibitors", *Nature*. **582**, 7811, pp 289-293
- [203] De P., Bhayye S., Kumar V., and Roy K. (2020) "In silico modeling for quick prediction of inhibitory activity against 3CLpro enzyme in SARS CoV diseases", *Journal of Biomolecular Structure and Dynamics*. pp 1-27
- [204] Toropov A.A., Toropova A.P., Veselinović A.M., Leszczynska D., et al. (2020) "SARS-CoV Mpro inhibitory activity of aromatic disulfide compounds: QSAR model", *Journal of Biomolecular Structure and Dynamics*. pp 1-7
- [205] Khaerunnisa S., Kurniawan H., Awaluddin R., Suhartati S., et al. (2020) "Potential inhibitor of COVID-19 main protease (Mpro) from several medicinal plant compounds by molecular docking study", *Preprints*. **2020**, pp 2020030226
- [206] Cao B., Wang Y., Wen D., Liu W., et al. (2020) "A trial of lopinavir–ritonavir in adults hospitalized with severe Covid-19", *New England Journal of Medicine*. pp 29-39
- [207] Abdelrheem D.A., Ahmed S.A., Abd El-Mageed H., Mohamed H.S., et al. (2020) "The inhibitory effect of some natural bioactive compounds against SARS-CoV-2 main protease: insights from molecular docking analysis and molecular dynamic simulation", *Journal of Environmental Science and Health, Part A*. **55**, 11, pp 1373-1386
- [208] Hasan A., Paray B.A., Hussain A., Qadir F.A., et al. (2021) "A review on the cleavage priming of the spike protein on coronavirus by angiotensin-converting enzyme-2 and furin", *Journal of Biomolecular Structure and Dynamics*. **39**, 8, pp 3025-3033
- [209] Murugan N.A., Kumar S., Jeyakanthan J., and Srivastava V. (2020) "Searching for target-specific and multi-targeting organics for Covid-19 in the Drugbank database with a double scoring approach", *Scientific reports*. **10**, 1, pp 1-16
- [210] Rehman M., AlAjmi M.F., and Hussain A. (2021) "Natural compounds as inhibitors of SARS-CoV-2 main protease (3CLpro): A molecular docking and simulation approach to combat COVID-19", *Current Pharmaceutical Design*. pp 10-19
- [211] Yu R., Chen L., Lan R., Shen R., et al. (2020) "Computational screening of antagonists against the SARS-CoV-2 (COVID-19) coronavirus by molecular docking", *International Journal of Antimicrobial Agents*. **56**, 2, pp 106012
- [212] Saíz-Urra L., González M.P., and Teijeira M. (2007) "2D-autocorrelation descriptors for predicting cytotoxicity of naphthoquinone ester derivatives against oral human epidermoid carcinoma", *Bioorganic & medicinal chemistry*. **15**, 10, pp 3565-3571
- [213] Sun M., Chen J., Cai J., Cao M., et al. (2010) "Simultaneously Optimized Support Vector Regression Combined With Genetic Algorithm for QSAR Analysis of KDR/VEGFR- 2 Inhibitors", *Chemical biology & drug design*. **75**, 5, pp 494-505
- [214] Jin X., Peldszus S., and Huck P.M. (2015) "Predicting the reaction rate constants of micropollutants with hydroxyl radicals in water using QSPR modeling", *Chemosphere*. **138**, pp 1-9
- [215] Kertesz D.J., Brotherton-Pleiss C., Yang M., Wang Z., et al. (2010) "Discovery of



piperidin-4-yl-aminopyrimidines as HIV-1 reverse transcriptase inhibitors. N-benzyl derivatives with broad potency against resistant mutant viruses", *Bioorganic & medicinal chemistry letters*. **20**, 14, pp 4215-4218

[216] Wong K.Y., Mercader A.G., Saavedra L.M., Honarparvar B., et al. (2014) "QSAR analysis on tacrine-related acetylcholinesterase inhibitors", *Journal of biomedical science*. **21**, 1, pp 1-8

[217] Chiari L.P., da Silva A.P., de Oliveira A.A., Lipinski C.F., et al. (2021) "Drug design of new sigma-1 antagonists against neuropathic pain: A QSAR study using partial least squares and artificial neural networks", *Journal of Molecular Structure*. **1223**, pp 129156

[218] Mandelker D., Gabelli S.B., Schmidt-Kittler O., Zhu J., et al. (2009) "A frequent kinase domain mutation that changes the interaction between PI3K $\alpha$  and the membrane", *Proceedings of the National Academy of Sciences*. **106**, 40, pp 16996-17001

[219] Zefirov N.S., and Palyulin V.A. (2001) "QSAR for boiling points of "small" sulfides. Are the "high-quality structure-property-activity regressions" the real high quality QSAR models?", *Journal of chemical information and computer sciences*. **41**, 4, pp 1022-1027

[220] Pourbasheer E., Beheshti A., Khajehsharifi H., Ganjali M.R., et al. (2013) "QSAR study on hERG inhibitory effect of kappa opioid receptor antagonists by linear and non-linear methods", *Medicinal Chemistry Research*. **22**, 9, pp 4047-4058

[221] Goudarzi N., and Goodarzi M. (2010) "Application of successive projections algorithm (SPA) as a variable selection in a QSPR study to predict the octanol/water partition coefficients ( $K_{ow}$ ) of some halogenated organic compounds", *Analytical Methods*. **2**, 6, pp 758-764

[222] Schuur J.H., Selzer P., and Gasteiger J. (1996) "The coding of the three-dimensional structure of molecules by molecular transforms and its application to structure-spectra correlations and studies of biological activity", *Journal of Chemical Information and Computer Sciences*. **36**, 2, pp 334-344

[223] Saíz-Urra L., González M.P., and Teijeira M. (2006) "QSAR studies about cytotoxicity of benzophenazines with dual inhibition toward both topoisomerases I and II: 3D-MoRSE descriptors and statistical considerations about variable selection", *Bioorganic & medicinal chemistry*. **14**, 21, pp 7347-7358

Lipinski's rule of five were calculated for all proposed compounds and the suggested new compounds have acceptable pharmacological properties. In the third study, the combination of the least absolute deviation-least absolute shrinkage and selection operator (LAD-LASSO) was introduced as a new variable selection method for the ANN-based QSAR studies. The ANN model coupled with the efficient LAD-LASSO variable selection method was evaluated by predicting the biological activity of three datasets of chemical compounds. The LAD-LASSO variable selection method was applied, and the descriptors with the most relevance to the biological activities were selected. The selected descriptors were defined as the ANN inputs, and the designed models were optimized. The biological activity of the test set compounds was predicted using the corresponding optimum ANN models. The coefficients of determination ( $R^2$ ) for the test data in the three datasets were equal to 0.87, 0.84, and 0.87. The applicability domain and Y-randomization test also proved the efficiency of the developed models. Finally, the created QSAR models were utilized to suggest numerous novel, highly potent chemical compounds by structurally modifying the weak molecules in all three datasets. The response value of the new suggested compounds was predicted using the optimum ANN models. According to the LR information, and the presence of different hydrophilic and hydrophobic interactions in the active site of the respective receptor indicates the high potential of chemical compounds. In the last work of this study, a combination of the SCAD and the ANN was utilized in the quantitative structure-retention indices(RIs) relationship (QSRR). The proposed SCAD method reduces the dimension of data before using the robust ANN modeling method. The efficiency of the SCAD-ANN methods was evaluated by the construction of a QSRR model between the most effective molecular descriptors and RIs for two sets of volatile organic compounds (VOCs). The SCAD method was applied to training data, and effective descriptors were selected in  $\lambda_{\min}$  and were defined as the inputs to the ANN modeling method. All ANN parameters were optimized simultaneously. Some statistical parameters were computed, and the obtained results indicate that the constructed QSRR models have acceptable values. Also, the applicability domain analysis reveals that more than 95% of the data are in the confidence range. So, the prediction results of the SCAD-ANN models are reliable.

**Keywords:** QSAR, QSPR, Cancer, HIV, Coronavirus, SCAD, ALASSO, LAD-LASSO, ANN, Molecular docking

## Abstract

The important goal of this dissertation is to use the penalized variable selection methods to select the most effective descriptors and coupling the penalized methods with the nonlinear artificial neural network (ANN) modeling. In this regard, different methods were used to make a relationship between the selected descriptors and the target response. So, in the first study, a hybrid of smoothly clipped absolute deviation (SCAD) and ANN was used as a new approach (SCAD-LM-ANN) in the quantitative structure-activity relationship (QSAR) studies. SCAD-LM-ANN exploits the useful shrinkage nature of penalized SCAD method in the reduction of high dimensional data prior to the modeling method. The performance of the SCAD-LM-ANN model was examined by establishing a QSAR model between Dragon derived descriptors and biological activities for a set of thioacetamide/acetanilide derivatives as HIV inhibitors. SCAD method with the 10-fold cross validation was applied to the dataset in the absence of test set compounds. A number of 11 descriptors were selected at a  $\lambda$  with the lowest cross-validation error ( $\lambda_{\min}$ ). The selected descriptors were used as inputs of the ANN. All parameters affecting model performance were optimized, and the SCAD-LM-ANN model with the architecture of 5-5-1 was selected as the optimal QSAR model. Several statistical parameters were considered for the model evaluation. The obtained results prove the generalizability and predictability of the proposed SCAD-LM-ANN model. According to the established relationship in the recommended QSAR model, new derivatives were designed and suggested as new active HIV inhibitors for further studies. The accuracy of the suggested compounds was studied and confirmed by analyzing the ligand-receptor (LR) interactions derived from the molecular docking studies. In the second part of this dissertation, a new approach as a hybrid of adaptive least absolute shrinkage and selection operator (ALASSO) and ANN was introduced for the 3-chymotrypsin-like protease (3CL<sup>Pro</sup>) inhibitors as SARS CoV-2 potent compounds. Given the importance of this disease since 2019, the recommendation and design of new active compounds are crucial. After evaluating the accuracy and validity of the developed ALASSO-LM-ANN model, new compounds were proposed using effective descriptors, and the biological activity of the new compounds was predicted. The investigation of LR interactions were also performed using molecular docking study. The PK properties and



Faculty of Chemistry

Ph.D. Thesis in Analytical Chemistry

**Application of penalized regression techniques as new methods of variable selection in the quantitative structure-activity (QSAR) and structure-property (QSPR) studies**

By:

**Zeinab Mozafari**

Supervisor:

**Dr. Mansour Arab Chamjangali**

Advisors:

**Dr. Mohammad Arashi**

**Dr. Nasser Goudarzi**

February 2022