



## مراحل، محاسبات و نتایج حاصل از مطالعات پیش‌بینی‌های نظری ارتباط کمی ساختار بازداری (QSRR) اسانس گیاه میخک زینتی

مجید محمدحسینی<sup>۱\*</sup>، مهدی نکوئی<sup>۱</sup>

<sup>۱</sup>گروه شیمی، واحد شاهرود، دانشگاه آزاد اسلامی، شاهرود، ایران

تاریخ ثبت اولیه: ۱۴۰۲/۰۱/۲۵، تاریخ دریافت نسخه اصلاح شده: ۱۴۰۲/۰۴/۲۲، تاریخ پذیرش قطعی: ۱۴۰۲/۰۴/۳۰

### چکیده

در این مقاله، به تشریح مبسوط مدل‌های خطی توانمند در پیش‌بینی شاخص بازداری کوآتس گروه وسیعی از ترکیبات طبیعی شناسایی شده در روغن اسانسی گیاه میخک زینتی به عنوان یکی از گیاهان دارویی پرداخته شده است. در این راستا، اساس کار مبتنی بر روابط کمی ساختار بازداری (QSRR) می‌باشد که در منابع علمی از اهمیت بسزایی جهت برقراری ارتباط منطقی و هدفمند بین شاخص کوآتس به عنوان یک متغیر وابسته و گروهی از توصیف‌کننده‌های مولکولی به عنوان متغیرهای مستقل برخوردار است. در این راستا، پس از ترسیم ساختار ترکیبات مفروض در محیط نرم‌افزار هایپرکم و بهینه‌سازی آن‌ها، جهت استخراج توصیف‌کننده‌های مولکولی مربوطه از نرم‌افزار دراگون استفاده شد. در مرحله بعد، پس از حذف توصیف‌کننده‌های غیر مرتبط و اضافی، نهایتاً با روش‌های مرحله‌ای و روش انتخاب متغیر مبتنی بر الگوریتم ژنتیک گروهی از توصیف‌کننده‌های مهم و مؤثر شناسایی و ارتباط خطی آن‌ها با شاخص بازداری کوآتس مورد بحث و بررسی قرار گرفت. نتایج حاصله حاکی از توانمندی بالای مدل‌های ارائه شده جهت پیش‌بینی شاخص کوآتس گروه وسیعی از ترکیبات طبیعی دارد.

واژه‌های کلیدی: ارتباط کمی ساختار- بازداری، رگرسیون خطی چندگانه، توصیف‌کننده‌های مولکولی، الگوریتم ژنتیک، شاخص بازداری کوآتس، میخک زینتی.

### ۱. مقدمه

هم‌زمان با پیدایش انسان‌ها، استفاده از گیاهان دارویی نیز آغاز شد. برخی از ترکیبات شیمیایی موجود در گیاه دارای ساختار پیچیده‌ای هستند و سنتز آن‌ها در آزمایشگاه یا غیرممکن یا با صرف زمان و هزینه زیاد امکان‌پذیر است. در قرن ۱۸ و اوایل قرن ۱۹، محققان پیشرفت قابل توجهی در خالص‌سازی و شناسایی ترکیبات شیمیایی موجود در گیاهان داشته و موادی را به صورت فرآورده-

\*مهم‌دار مکاتبات: مجید محمدحسینی

نشانی: گروه شیمی، واحد شاهرود، دانشگاه آزاد اسلامی، شاهرود، ایران

پست الکترونیک: majidmohammadhosseini@yahoo.com: E-mail

تلفن: ۰۲۳۲۲۳۹۴۲۷۸

های داروئی برای مصرف عرضه کردند. هم‌زمان با انقلاب صنعتی، علم شیمی پیشرفت چشمگیری داشت که باعث به‌وجود آمدن این تفکر در محیط‌های علمی شد که می‌توان از طریق سنتز ترکیبات شیمیایی به‌خصوص مواد داروئی مشکل دارو و درمان بیماری‌ها را حل کرد [۱]. به همین دلیل تولید داروهای شیمیایی در قرن بیستم سرعت روزافزونی پیدا کرد و داروهای گیاهی به‌دست فراموشی سپرده شدند. پس از مواجه شدن با مشکلاتی نظیر آلودگی آب و هوا و خاک که توسط کارخانجات تولید مواد شیمیایی ایجاد شده بوده و عوارض جانبی داروهای شیمیایی که بعضاً پس از چند نسل ظاهر می‌شوند، به فکر استفاده از مواد طبیعی فن-آوری‌های غیر مخرب افتادند. به‌طوری‌که در کشورهای صنعتی مصرف داروهای گیاهی از مرز ۷ درصد گذشته است.

با توجه به کاربرد داروهای گیاهی، پژوهشکده‌ی صنایع شیمیایی سازمان پژوهش‌های علمی و صنعتی ایران نیز تحقیق‌هایی را در این زمینه به اجرا در آورده است. تاریخچه‌ی گیاهان داروئی به طور دقیق مشخص نیست و در طول تاریخ استفاده از گیاهان داروئی با خرافات و آداب خاصی همراه بوده است. مصری‌ها و چینی‌ها از اولین اقوامی هستند که از حدود ۲۷۰۰ سال پیش از میلاد مسیح از گیاهان به عنوان دارو استفاده می‌کردند. تئوفراست یکی از شاگردان ارسطو بنیانگذار مکتب درمان با گیاه است. دیوسکورید در قرن اول میلادی مجموعه‌ای را مشتمل بر خواص داروئی ۶۰۰ گیاه جمع‌آوری نمود که این اثر منشأ بسیاری از مطالعات در قرون بعد گردید [۱]. در قرون هشتم تا دهم میلادی، بوعلی‌سینا و محمد زکریای رازی سبب توسعه‌ی دانش درمان با گیاه شدند و در قرن سیزدهم، ابن بیطار خصوصیات بیش از ۱۴۰۰ گیاه را در کتابی گردآوری نموده و در قرن نوزدهم داروهای شیمیایی به سرعت جایگزین بسیاری داروهای گیاهی گردید. سپس، در اواخر قرن بیستم عوارض جانبی و مضر داروهای شیمیایی سبب رویکرد دوباره دانشمندان به گیاهان داروئی شد به‌طوری‌که این دوره را رنسانس گیاهان داروئی نامیدند. تا قرن نوزدهم، گیاهان داروئی به شکل بسیار ابتدایی مورد مصرف قرار می‌گرفتند تا این‌که استخراج مواد مؤثر گیاهی، از قرن نوزدهم آغاز گردید [۱].

نظر به کاربرد گیاهان داروئی در زمینه‌های مختلف، گیاهان اساساً در سه گروه اصلی گیاهان داروئی، گیاهان ادویه‌ای، گیاهان عطری طبقه‌بندی می‌شوند. اما به‌طور کلی، گیاه داروئی به گیاهی اطلاق می‌شود که در یک یا چند اندام خود حاوی مواد مؤثر بوده و کاشت، داشت و برداشت آن‌ها صرفاً به منظور استفاده از این مواد انجام گردد. بنابراین، با آن‌که اندام‌های برخی گیاهان نظیر برگ‌های گردو و کاکل ذرت و پوست میوه لوبیا حاوی مواد مؤثری هستند که کاربردهای داروئی نیز دارند ولی از این نظر که کاشت، داشت و برداشت این گیاهان تنها به منظور استفاده از این گیاهان تنها به منظور استفاده از مواد مؤثر موجود در آن‌ها انجام نمی‌گیرد، گیاه داروئی محسوب نمی‌شوند. نظر به تنوع و اهمیت حیاتی گیاهان داروئی، یکی از جامع‌ترین دایره‌المعارف‌های گیاهی توسط دکتر ولی‌الله مظفریان به زیور طبع آراسته شده است [۲].

یکی از مهم‌ترین مواد مؤثر در گیاهان داروئی را روغن‌های فرار یا روغن‌های اسانسی تشکیل می‌دهند. این مواد در قسمت‌های مختلف بسیاری از گیاهان داروئی وجود دارند. بسیاری از گیاهان داروئی به‌علت داشتن روغن‌های فرار به‌طور مستقیم در پزشکی مصرف می‌شوند، ولی در بیشتر موارد روغن‌های فرار را از مواد خام جدا نموده و به‌عنوان دارو به کار می‌برند. معمولاً، از اسانس‌ها

در تهیه‌ی عطرها و اسپری‌های خوشبوکننده و هم‌چنین به عنوان معطرکننده‌ی صابون و خمیردندان و شامپوهای طبی مورد استفاده قرار می‌گیرند. از اسانس‌ها در صنایع داروسازی نیز استفاده‌های زیادی به‌عمل می‌آید، زیرا بسیاری از این مواد دارای اثر ضد باکتری و ضد عفونی‌کننده هستند. از اثرهای داروئی اسانس‌ها، می‌توان به خواصی مانند ضد تورم، ضد دل‌درد، آرام‌بخش، ضد نفخ، اشتها آور و خلط‌آور اشاره نمود. در صنایع غذایی و کنسروسازی، روغن‌های اسانسی به‌دست آمده از گیاهان ادویه‌ای برای بهبود طعم مواد غذایی استفاده می‌شوند. هم‌چنین، روغن‌های اسانسی می‌توانند ضمن بهبود طعم برخی از داروها به‌عنوان محافظ هم مورد استفاده قرار گیرند.

هدف از تحقیق اخیر، مدل‌سازی‌های توانمند مبتنی بر روابط کمی ساختار-بازداری (QSRR) ۶۰ ترکیب طبیعی شناسایی شده در روغن اسانسی گیاه *Pittosporum undulatum* می‌باشد. این مدل‌سازی، با استفاده از روش‌های رگرسیون خطی چندگانه مرحله‌ای (SW-MLR) و الگوریتم ژنتیک-رگرسیون خطی چندگانه (GA-MLR) انجام شده است.

## ۲. روش‌های محاسباتی

### ۱-۲. سری داده‌ها

در این مطالعه، اندیس کواتس ۶۰ ترکیب از اسانس گیاه داروئی میخک زینتی (شکل ۱)، با استفاده از روش رگرسیون خطی چندگانه مورد بررسی قرار گرفته است. سری داده‌ی مربوطه، شامل مجموعه‌ای از ترکیبات مونوترپن‌ها، سزکوئی‌ترین‌ها و سایر ترکیبات طبیعی، از تحقیق یک گروه از محققین پرتقالی به سال ۲۰۰۷ میلادی گرفته شده است [۳]. مراحل انجام این پژوهش، شامل انتخاب سری داده‌ها، تعیین توصیف‌کننده‌ها، رگرسیون خطی چندگانه و مدل‌سازی می‌باشد. نام، ساختار و مقادیر عددی اندیس‌های بازداری (RI) آن‌ها در جدول ۱، نشان داده شده است. برای مدل‌سازی، ابتدا اندیس بازداری ترکیبات، به‌عنوان پارامتر مورد بحث، انتخاب شد. در مرحله‌ی دوم، ساختار تمام گونه‌های مورد نظر رسم و بهینه گردید.



شکل ۱. تصویر گیاه میخک زینتی

## ۲-۲. محاسبه و انتخاب توصیف کننده‌ها

توصیف کننده‌ها، کدهای عددی هستند که بیان گر ویژگی‌های خاصی از هر مولکول بوده که با خواص فیزیکی-شیمیایی ماده - ارتباط دارند. با توجه به موضوع مورد بحث (بررسی ارتباط کمی ساختار- بازداري)، ابتدا ساختار مولکول‌ها به وسیله نرم افزار Hyperchem رسم و ساختار هندسی آن‌ها با احتساب اتم‌های هیدروژن و با استفاده از مدل هامیلتونی AM1 و الگوریتم Polak-Ribiere که از روش‌های نیمه تجربی می‌باشند، بهینه شد. سپس، ساختارهای رسم شده، توسط نرم افزار مذکور به عنوان ورودی به نرم افزار دیگری به نام Dragon، معرفی شد. محاسبه‌ی توصیف کننده‌ها، توسط این نرم افزار صورت گرفت. بالغ بر ۱۴۰۰ توصیف کننده از هجده گروه مختلف برای هر ترکیب با این بسته‌ی نرم افزاری، قابل محاسبه است.

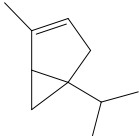
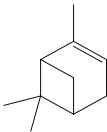
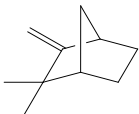
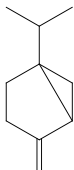
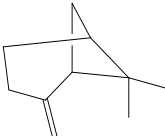
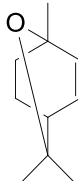
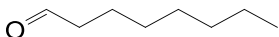
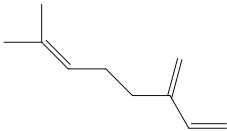
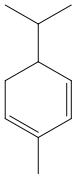
## ۲-۳. تجزیه و تحلیل آماری توصیف کننده‌ها

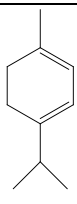
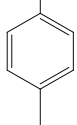
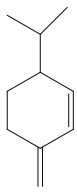
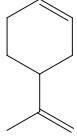
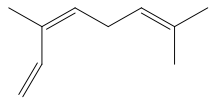
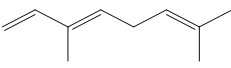
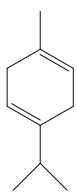
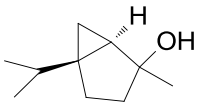
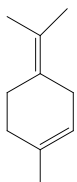
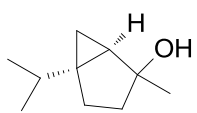
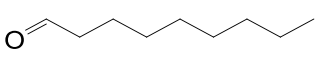
برای ایجاد مدلی که بیان گر ارتباط ساختاری ترکیبات مورد بررسی با اندیس بازداري آن‌ها باشد، از روش رگرسیون خطی چندگانه (MLR) استفاده شد. برای مشاهده و مرتب کردن توصیف کننده‌ها، خروجی نرم افزار Dragon را به نرم افزار Excel انتقال می‌دهیم. طبیعی است که تعداد زیاد توصیف کننده‌ها باعث پیچیدگی محاسبات شده و هم چنین احتمال وجود فاکتورهای دارای برهم کنش باهم را افزایش می‌دهد. لذا، تعدادی از این توصیف کننده‌ها که تمام مقادیر آن‌ها صفر و یا بالای ۹۰٪ مقادیر یکسان بودند حذف شدند. نهایتاً، پس از مرحله‌ی حذف مقدماتی اولیه، توصیف کننده‌های باقیمانده در یک ماتریس داده‌ی  $n \times m$  جای گرفتند. در این ماتریس، جملات  $n$  و  $m$ ، به ترتیب بیان گر تعداد ترکیبات و تعداد توصیف کننده‌ها هستند. مقادیر عددی این دو جمله به ترتیب برابر ۶۰ و ۱۹۸ و شکل ماتریس مربوطه به صورت  $60 \times 198$  می‌باشد. در پایان، بردار ستونی  $(y)$  برای محاسبه‌ی متغیر وابسته (اندیس کواتس) ساخته شد.

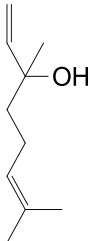

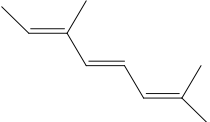
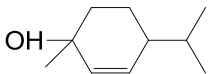
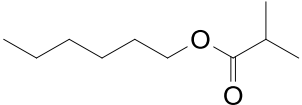
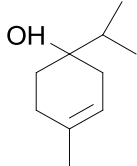
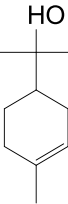
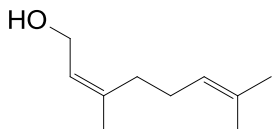
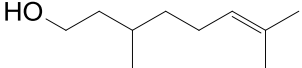
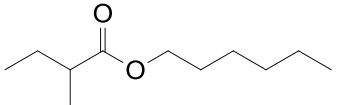
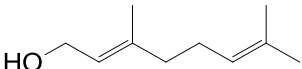
## ۲-۴. الگوی پراکنش آنالیز جزء اصلی (PCA) و تقسیم سری داده

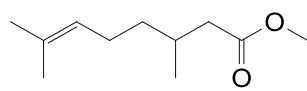
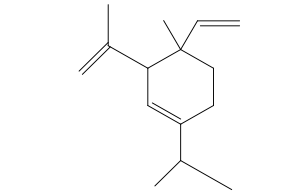
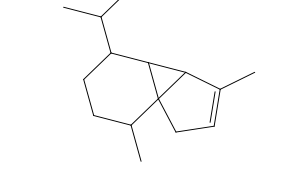
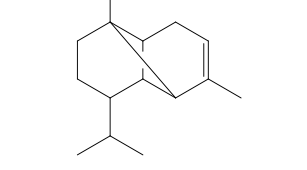
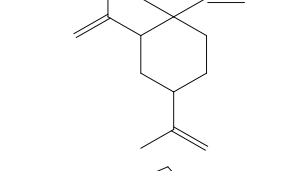
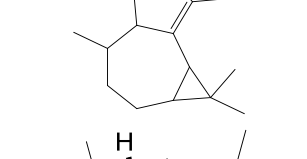
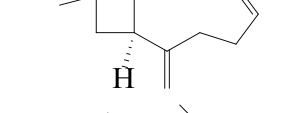
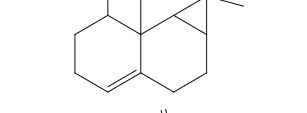
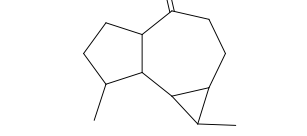
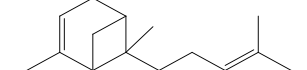
در حقیقت، PCA یک راهکار مطمئن در تجزیه و تحلیل و طبقه‌بندی اطلاعات محسوب می‌شود. PCA، می‌تواند چند جزء زمینه-ای اصلی را تعیین نموده و به تشریح واریانس مسبوط مشاهده شده در اکثر داده‌ها کمک شایانی نماید [۴،۵]. هدف اصلی در هر آنالیز مبتنی بر PCA، مشخص ساختن هر شیء در ماتریس ورودی (ردیف‌ها) بدون تجزیه‌ی هر گونه متغیر (ستون‌ها) می‌باشد. در این راستا، داده‌ها در یک زیر مجموعه‌ی بسیار کوچک تر از متغیرهای جدید و یا امتیازات مؤلفه‌های اصلی پیش‌بینی می‌شوند. این متغیرها یا فاکتورهای جدید، ارتباطات خطی مقادیر اولیه بوده و بیان گر واریانس مربوط به هر سری داده هستند که می‌توانند موارد اضافی را حذف نمایند. در حقیقت، اجزاء اصلی متوالی، در مرتبه‌ی کاهشی مقادیر ویژه مرتب می‌شوند [۶]. اجزاء اصلی یا PCها قائم یا مستقل بوده و به نحوی مقیاس بندی می‌شوند که واریانس‌های مربوطه با مقدار واحد برابر شوند. هم چنین، اجزاء اصلی به گونه‌ای آرایش می‌یابند که واریانس بیان شده با اولین PC بیشترین باشد و واریانس‌های مربوط به PC دوم، سوم و ... به ترتیب روند کاهشی داشته باشند. بدین ترتیب، آخرین واریانس کم ترین مقدار را خواهد داشت [۷]. در ادامه پس از اصلاح و کاهش توصیف کننده‌ها، با اعمال دستور متنی آنالیز جزء اصلی (PCA) در محیط Matlab، شکل (۲) حاصل شد. در این شکل، هر علامت + نشان گر یک ترکیب و شماره‌ی مندرج در کنار آن، بیان کننده‌ی شماره‌ی آن ترکیب است.

جدول ۱. نام، ساختار و اندیس کواتس ترکیبات اسانس گیاه میخک زینتی

شماره	نام ترکیب	RI	ساختار ترکیب
۱	$\alpha$ -Thujene	924	
۲	$\alpha$ -Pinene	930	
۳	Camphene	938	
۴	Sabinene	958	
۵	$\beta$ -Pinene	963	
۶	Dehydro-1,8-Cineole	973	
۷	<i>n</i> -Octanal	973	
۸	Myrcene	975	
۹	$\alpha$ -Phellandrene	995	

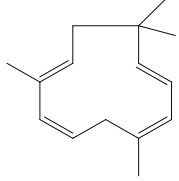
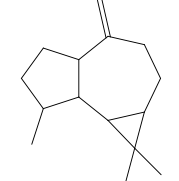
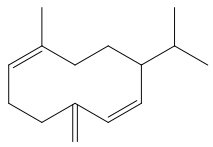
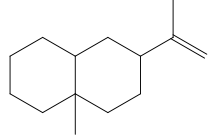
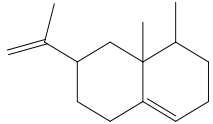
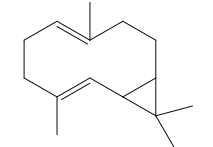
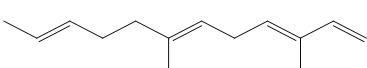
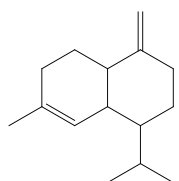
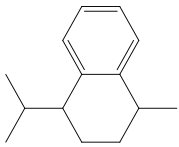
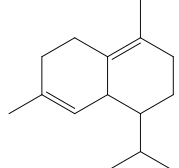
	1002	$\alpha$ -Terpinene	۱۰
	1003	<i>p</i> -Cymene	۱۱
	1005	$\beta$ -Phellandrene	۱۲
	1009	Limonene	۱۳
	1017	<i>cis</i> - $\beta$ -Ocimene	۱۴
	1027	<i>trans</i> - $\beta$ -Ocimene	۱۵
	1035	$\gamma$ -Terpinene	۱۶
	1037	<i>trans</i> -Sabinene hydrate	۱۷
	1064	Terpinolene	۱۸
	1066	<i>cis</i> -Sabinene hydrate	۱۹
	1073	<i>n</i> -Nonanal	۲۰

	1074	Linalool	۲۱
	1095	<i>trans-p</i> -Menth-2-en-1-ol	۲۲
	1110	<i>allo</i> -Ocimene	۲۳
	1110	<i>cis-p</i> -Menthen-1-ol	۲۴
	1127	Hexyl isobutyrate	۲۵
	1148	Terpinene-4-ol	۲۶
	1159	$\alpha$ -Terpineol	۲۷
	1206	Nerol	۲۸
	1208	Citronellol	۲۹
	1222	Hexyl-2-Methylbutyrate	۳۰
	1236	Geraniol	۳۱

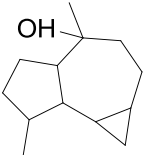
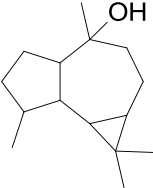
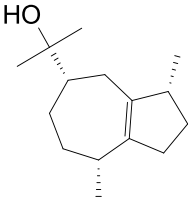
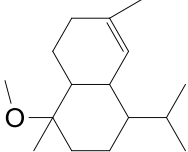
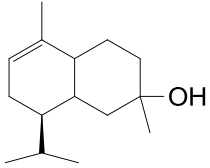
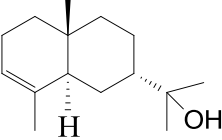
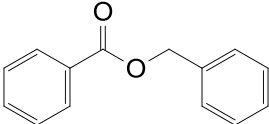
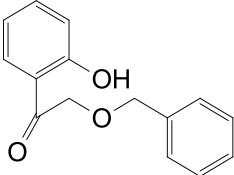
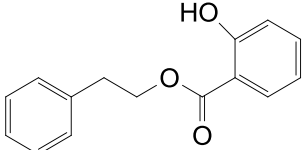
	1245	Methyl citronellate	۳۲
	1332	$\delta$ -Elemene	۳۳
	1345	$\alpha$ -Cubebene	۳۴
	1375	$\alpha$ -Copaene	۳۵
	1388	$\beta$ -Elemene	۳۶
	1400	$\alpha$ -Gurjunene	۳۷
	1414	$\beta$ -Caryophyllene	۳۸
	1426	$\beta$ -Gurjunene	۳۹
	1428	Aromadendrene	۴۰
	1434	<i>trans</i> - $\alpha$ -Bergamotene	۴۱



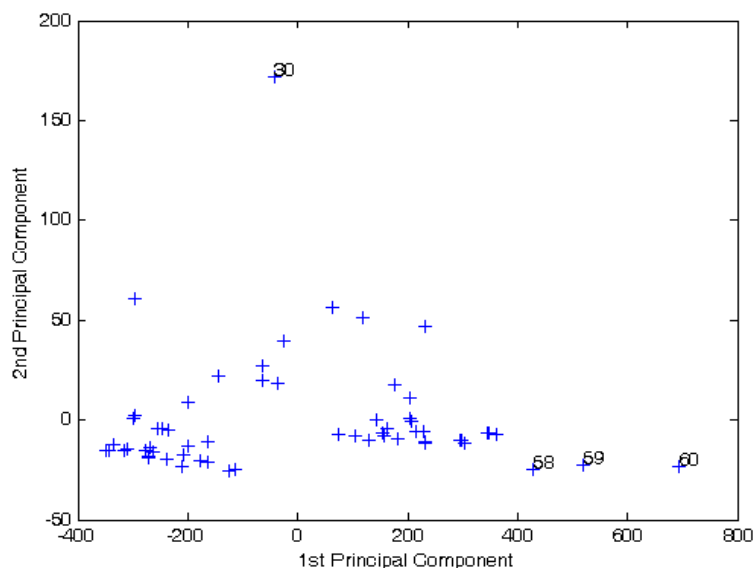
---

	1447	$\alpha$ -Humulene	۴۲
	1454	<i>allo</i> -Aromadendrene	۴۳
	1474	Germacrene-D	۴۴
	1476	$\beta$ -Selinene	۴۵
	1477	Valencene	۴۶
	1487	Bicyclogermacrene	۴۷
	1500	$\alpha$ - <i>trans-trans</i> -Farnesene	۴۸
	1500	$\gamma$ -Cadinene	۴۹
	1505	Calamenene	۵۰
	1505	$\delta$ -Cadinene	۵۱

---

	1566	Globulol	۵۲
	1569	Viridiflorol	۵۳
	1575	Guaiol	۵۴
	1616	T-Cadinol	۵۵
	1618	$\delta$ -Cadinol	۵۶
	1634	$\alpha$ -Eudesmol	۵۷
	1701	Benzyl benzoate	۵۸
	1790	Benzyl salicylate	۵۹
	1966	Phenylethyl salicylate	۶۰

بر اساس این شکل، قسمت عمده‌ی ترکیبات در خوشه‌ی اصلی جای گرفته ولی ترکیبات دارای شماره‌های ۳۰، ۵۸، ۵۹ و ۶۰ دارای فاصله‌ی محسوس با سایر ترکیبات در سری داده‌ی مفروض هستند. این شکل، اساساً دارای دو پیام است: نخست اینکه احتمال دارد این ترکیبات دارای رفتاری متفاوت با ترکیبات دیگر باشند که این امر در برخی از موارد به حصول نتایج نامناسب و عدم مدل‌سازی مناسب منجر می‌شود. دوم اینکه ۴ ترکیب فوق را نمی‌توان در سری آزمون در نظر گرفت، چرا که در سری آزمون باید ترکیبات دارای رفتار و ماهیت مشابه یا نزدیک با ترکیبات موجود در خوشه‌ی اصلی قرار گیرند. نهایتاً، با در نظر گرفتن الگوی پراکنش نقطه-ای آنالیز مؤلفه‌ی اصلی، جهت تقسیم هوشمند سری داده به مجموعه‌های آموزش و آزمون ترکیباتی در سری آموزش جای گرفتند که نماینده‌ی مناسبی از کل ترکیبات در سری داده‌ی اصلی باشند. وظیفه‌ی خطیر مدل‌سازی، بر دوش ترکیبات سری آموزش قرار داشته و ترکیبات سری آزمون در این امر دخالت ندارند. در عوض، قدرت پیشگویی مدل و توانمندی آن با اعمال به سری آزمون مورد بررسی قرار می‌گیرد. در مدل‌سازی مبتنی بر رگرسیون خطی چندگانه روی اسانس گیاه دارویی میخک زینتی، سری اطلاعات اولیه‌ی مشکل از ۶۰ مولکول را به مجموعه‌های آموزشی شامل ۵۱ ترکیب و آزمایشی شامل ۹ ترکیب تقسیم می‌کنیم. این تقسیم-بندی با نسبت تقریبی ۴ به ۱ صورت می‌پذیرد که در نهایت ۸۵٪ از سری اطلاعات در مجموعه‌ی آموزشی و ۱۵٪ از آن‌ها در مجموعه‌ی آزمایشی قرار می‌گیرند.



شکل ۲. الگوی پراکنش نقطه‌ای ترکیبات سری داده‌ی اصلی میخک زینتی با استفاده از دستور آنالیز جزء اصلی (PCA)

### ۳. نتایج و بحث

#### ۳-۱. تحلیل مدل‌های آماری و انتخاب مدل مناسب

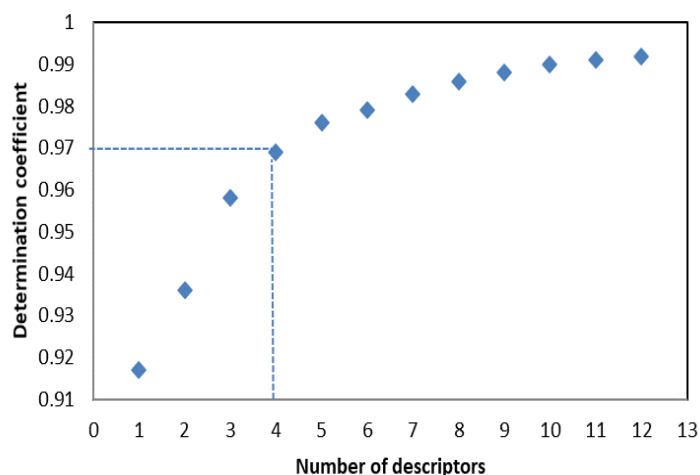
برای به‌دست آوردن مدلی مناسب، سری اطلاعات مجموعه‌ی آموزشی را به‌عنوان متغیر مستقل و مقادیر تجربی اندیس بازداري را به‌عنوان متغیر وابسته به نرم‌افزار آماری SPSS معرفی کرده و با استفاده از منوی آنالیز، گزینه‌ی رگرسیون خطی با روش مرحله‌ای مدل‌سازی صورت می‌پذیرد. نهایتاً، چندین مدل مختلف به‌طور جداگانه به‌دست می‌آید که با توجه به خصوصیات آماری آن‌ها، از جمله ضریب رگرسیون یا ضریب هم‌بستگی ( $R$ )، آماره‌ی  $F$  و خطای استاندارد (SE)، پس از رسم مقادیر  $R$  و  $R^2$  و  $SE$  برحسب

تعداد توصیف‌کننده‌ها بهترین مدل که دارای بیش‌ترین مقدار  $R$  و  $F$ ، و کم‌ترین مقدار خطای استاندارد و شامل توصیف‌کننده‌های تا حد امکان قابل توجیه باشد، به عنوان مدل نهایی برای ارتباط اندیس بازداری مولکول‌ها با ساختار آن‌ها انتخاب شد. بر اساس محاسبات نهایی، مجموعه‌ای شامل ۱۴ مدل رگرسیون در نرم‌افزار SPSS حاصل شد. نتایج عددی خلاصه‌ی مدل‌های ارائه شده در محیط SPSS، در جدول ۲ نمایش داده شده است.

در شکل ۳، روند تغییرات مقادیر ضرایب تعیین ( $R^2$ ) بر حسب تعداد توصیف‌کننده‌ها نشان داده شده است. به عنوان یک معیار کلی، در روش‌های رگرسیون خطی چندگانه، با رسم روند تغییر ضریب اندازه‌گیری به‌عنوان تابعی از توصیف‌کننده و تعیین نقطه‌ی شکست نمودار، می‌توان به تعداد توصیف‌کننده‌ی بهینه پی برد. بدیهی است هرچه تعداد توصیف‌کننده‌های انتخابی کم‌تر باشد، مدل انتخابی ساده‌تر بوده و از اعتبار بیشتری برخوردار است.

جدول ۲. خلاصه‌ی ۱۲ مدل منتخب پس از محاسبات رگرسیون در محیط نرم‌افزار SPSS سری داده‌ی اسانس گیاه میخک زینتی

شماره‌ی مدل	ضریب هم-بستگی ( $r$ )	ضریب اندازه‌گیری ( $r^2$ )	ضریب اندازه‌گیری تنظیم شده ( $Adj.r^2$ )	خطای استاندارد تخمین
۱	0.958	0.917	0.916	75.96
۲	0.967	0.936	0.933	67.75
۳	0.979	0.958	0.955	55.28
۴	0.985	0.969	0.967	47.84
۵	0.988	0.976	0.973	43.1
۶	0.990	0.979	0.977	40.04
۷	0.991	0.983	0.980	37.3
۸	0.993	0.986	0.984	33.4
۹	0.994	0.988	0.986	31.1
۱۰	0.995	0.990	0.988	29.1
۱۱	0.996	0.991	0.989	27.7
۱۲	0.996	0.992	0.990	26.1



شکل ۳. منحنی تغییرات مقادیر ضرایب تعیین ( $R^2$ ) بر حسب تعداد توصیف‌کننده‌ها با استفاده از نرم‌افزار SPSS سری داده‌ی اسانس گیاه میخک زینتی

در این مرحله، مدل انتخابی را برای مجموعه‌ی آزمایشی نیز اعمال می‌کنیم. از آنجایی که توصیف‌کننده‌ها باید متغیرهایی مستقل باشند، لذا برای اطمینان از عدم هم‌بستگی بین آن‌ها، در نرم‌افزار SPSS، با استفاده از منوی آنالیز، گزینه‌ی هم‌بستگی دو متغیره را انتخاب می‌کنیم. چنانچه هم‌بستگی بین توصیف‌کننده‌ها بیش‌تر از ۰/۸۵ باشد، توصیف‌کننده‌های وابسته را از سری اطلاعات حذف کرده و مراحل مذکور را تا جایی که هم‌بستگی بین توصیف‌کننده‌ها کم‌تر از این مقدار باشد، مجدداً انجام می‌دهیم. مقادیر عددی توصیف‌کننده‌های وارد شده در مدل، در جدول ۳ آورده شده است.

جدول ۳. مقادیر عددی توصیف‌کننده‌های وارد شده در مدل خطی نهایی سری داده‌ی اسانس گیاه میخک زینتی

شماره ترکیبات	توصیف‌کننده‌های نهایی انتخاب شده			
	LPRS	Mor29m	R7m+	MATS1e
Training set				
1	30.437	-0.208	0.001	0.059
2	30.023	-0.108	0	0.059
3	29.615	-0.213	0	0.059
4	30.437	-0.241	0.001	0.059
5	34.504	-0.198	0	0.048
6	29.343	-0.042	0.017	-0.004
7	33.542	-0.079	0.025	-0.064
8	31.601	-0.076	0.01	0
9	31.601	-0.08	0.009	0
10	31.601	0.009	0.01	0
11	31.601	-0.073	0.01	0
12	33.542	-0.064	0.029	-0.064
13	33.542	-0.044	0.021	-0.064
14	31.601	-0.019	0.009	0
15	51.62	-0.283	0.051	0.016
16	31.601	-0.077	0.009	0
17	35.153	-0.191	0.001	-0.066
18	34.733	-0.038	0.012	-0.004
19	36.453	-0.061	0.015	-0.072
20	33.542	0.008	0.019	-0.064

21	36.453	-0.079	0.014	-0.079
22	45.177	-0.095	0.02	0.017
23	36.011	-0.052	0.009	-0.072
24	36.32	-0.087	0.012	-0.072
25	38.959	0.054	0.019	-0.078
26	38.959	-0.024	0.018	-0.073
27	50.278	-0.179	0.024	0.019
28	38.959	0.054	0.019	-0.078
29	57.073	-0.173	0.033	0
30	55.962	-0.333	0.01	0.078
31	57.073	-0.12	0.029	0
32	55.604	-0.31	0.007	0.078
33	56.902	-0.325	0.013	0.04
34	55.604	-0.298	0.007	0.078
35	58.548	-0.095	0.033	0.04
36	57.938	-0.181	0.017	0
37	55.604	-0.425	0.009	0.078
38	58.08	-0.132	0.02	0
39	56.501	-0.071	0.02	0.04
40	57.423	-0.236	0.027	0.04
41	62.188	-0.037	0.021	-0.042
42	56.588	0.047	0.014	0.039
43	56.588	-0.147	0.014	0.04
44	60.845	-0.336	0.01	-0.049
45	60.845	-0.5	0.011	-0.049
46	61.991	-0.183	0.024	-0.054
47	61.818	-0.002	0.015	-0.054
48	62.168	-0.016	0.028	-0.054
49	65.768	0.12	0.032	0.001
50	71.272	0.132	0.046	-0.069

51	78.125	0.213	0.035	-0.06
Test set				
1	30.023	-0.077	0	0.059
2	31.601	-0.068	0.011	0
3	38.152	-0.099	0.025	-0.078
4	49.799	-0.017	0.017	0.017
5	55.742	-0.034	0.008	0.078
6	55.264	-0.265	0.006	0.078
7	56.677	-0.116	0.022	0.04
8	56.588	-0.07	0.014	0.04
9	61.818	-0.045	0.016	-0.054

## ۲-۳. حوزه‌ی کاربرد (AD)

در تحقیق اخیر، نمودار ویلیام یا نمودار روند تغییرات باقیمانده‌های استاندارد شده بر حسب لورج یا قدرت نفوذ [۸] جهت نمایش حوزه‌ی کاربرد فرآیند مدل‌سازی استفاده و نتایج حاصله در شکل ۴، نمایش داده شده است. طبق تعریف، لورج یا قدرت نفوذ نشان‌گر فاصله‌ی یک ترکیب از مرکز ثقل محور  $X$  می‌باشد. طبق یک قرارداد آماری، لورج یک ترکیب در فضای متغیر اولیه به صورت رابطه‌ی ۱، تعریف می‌شود [۹].

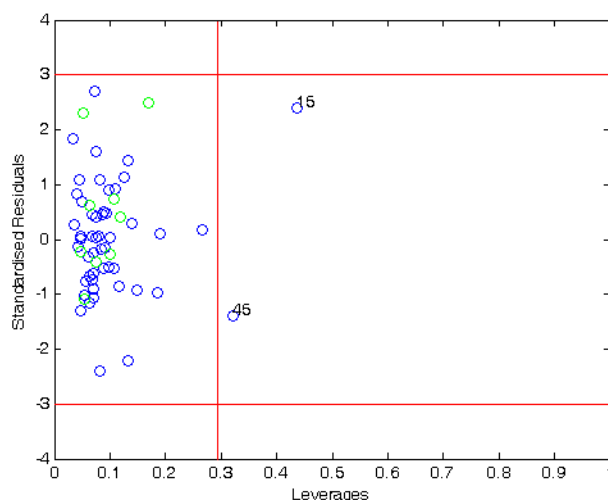
$$h_i = x_i^T (X^T X)^{-1} x_i \quad (1)$$

در این رابطه،  $x_i$  بردار توصیف‌کننده‌ی ترکیب مفروض و  $X$  ماتریس توصیف‌کننده‌ی حاصل از مقادیر توصیف‌کننده‌ی سری آموزش می‌باشد. مقدار لورج هشدار ( $h^*$ )، به صورت رابطه‌ی ۲ تعریف می‌شود [۱۰].

$$h^* = 3(p+1)/n \quad (2)$$

در رابطه‌ی (۲)،  $n$  تعداد ترکیبات در سری آموزش،  $p$  تعداد متغیرهای پیشگویی‌کننده می‌باشد. منظور از تعداد متغیرهای پیشگویی‌کننده، همان تعداد توصیف‌کننده‌های مولکولی منتخب در مدل‌سازی نهائی است. در فرآیندهای مدل‌سازی، ترکیبات دارای لورج بیشتر از لورج هشدار ( $h_i > h^*$ ) می‌توانند در بسیاری از موارد به شدت عملکرد مدل‌سازی را تحت تأثیر قرار دهند. در عین حال، در کلیه موارد نمی‌توان آن را به عنوان یک ترکیب انحرافی در نظر گرفت چرا که ممکن است باقیمانده‌ی استاندارد شده‌ی آن جزئی و قابل اغماض باشد. در شرایطی که برای یک ترکیب معین هم شرط  $h_i > h^*$  برقرار و هم باقیمانده‌ی استاندارد شده‌ی ترکیب مداری بزرگ و قابل توجه باشد، ترکیب مربوطه از حوزه‌ی کاربرد و در نتیجه سری داده‌ی اولیه حذف و مدل‌سازی با ترکیبات

باقیمانده انجام خواهد شد. تحت این شرایط، باید توجهی قانع کننده پیرامون علت حذف این ترکیب با در نظر گرفتن ساختار آن ارائه نمود. بعلاوه، مقدار ۳ برابر باقیمانده‌ی استاندارد شده معمولاً تحت عنوان مقدار برش برای قبول پیشگویی‌ها استفاده می‌شود. چرا که ثابت شده که نقاطی که در بازه‌ی  $\pm 3$  برابر باقیمانده‌های استاندارد شده از مقدار میانگین قرار می‌گیرند، ۹۹٪ اطلاعات دارای پراکنش نرمال را در برمی‌گیرند [۱۱]. بدین ترتیب، ترکیب لورج و باقیمانده‌های استاندارد شده برای مشخص نمودن حوزه‌ی کاربرد مورد استفاده قرار گرفت.



شکل ۴. نمودار ویلیام برای ترکیبات سری‌های آموزش و آزمون سری داده‌ی اسانس گیاه میخک زبنتی

بررسی شکل ۴، خاطر نشان می‌سازد که دو ترکیب (ترکیبات شماره‌ی ۱۵ و ۴۵) در سری آموزش دارای مقادیر لورج بیشتر از مقدار هشدار ( $h^*$ ) هستند. براساس رابطه‌ی ۲، مقدار لورج هشدار برابر با  $h^* = 0.2941$  است. بنابراین، می‌توانند جزو ترکیبات انحرافی باشند. خوشبختانه، در دو مورد اخیر مقادیر اطلاعات عددی پیشگویی شده مناسب و قابل قبول هستند. لذا، این دو ترکیب را لورج خوب خوانده و از حوزه‌ی کاربرد و سری داده‌ی اولیه حذف نمی‌کنیم. در یک جمع‌بندی، می‌توان گفت که نمودار حوزه‌ی کاربرد (نمودار ویلیام یا نمودار AD) هم تأییدکننده‌ی مناسب بودن مدل ساخته‌شده و هم بیان‌گر تقسیم مناسب سری داده‌ی اولیه به مجموعه‌های آموزش و آزمون است.

### ۳-۳. فاکتور تورم تغییر (VIF)

هم‌بستگی خطی چندگانه‌ی بین توصیف‌کننده‌های منتخب، با محاسبه‌ی پارامتری موسوم به فاکتور تورم تغییر (VIF) تعیین شد که در قالب رابطه‌ی ۳ قابل بیان است.

$$VIF = \frac{1}{1-r^2} \quad (3)$$



در رابطه‌ی اخیر،  $r$  ضریب هم‌بستگی رگرسیون چندگانه‌ی بین متغیرها در مدل نامیده می‌شود. طبق قرارداد، اگر VIF برابر یک باشد، هیچ‌گونه هم‌بستگی داخلی بین متغیرها وجود ندارد و در صورتی که بازه‌ی تغییرات VIF بین ۱ تا ۵ باشد، مدل مربوطه قابل قبول بوده و اگر مقادیر عددی محاسبه‌شده برای VIF بزرگ‌تر از ۱۰ باشد، مدل مربوطه ناپایدار بوده و یک بازنگری و تصحیح اساسی جهت توسعه‌ی آن ضروری خواهد بود [۱۲]. مقادیر VIF در جدول ۴ گزارش شده است.

### ۳-۴. اثر متوسط ( $MF_j$ ) توصیف‌کننده‌های منتخب

جهت بررسی اهمیت نسبی و تأثیر هر توصیف‌کننده در مدل، مقدار اثر متوسط ( $MF_j$ ) برای هر توصیف‌کننده محاسبه شد. محاسبه‌ی مقادیر عددی این پارامتر مهم آماری با استناد به رابطه‌ی ۴، صورت گرفت [۱۳].

$$MF_j = \frac{\beta_j \sum_{i=1}^n d_{ij}}{\sum_i \beta_j \sum_i d_{ij}} \quad (4)$$

در این رابطه،  $MF_j$  اثر متوسط توصیف‌کننده‌ی مفروض  $Z_j$ ،  $\beta_j$  ضریب توصیف‌کننده‌ی  $Z_j$  مقدار توصیف‌کننده‌های هدف برای هر مولکول و  $m$  شماره‌ی توصیف‌کننده در مدل ارائه شده می‌باشند. در حقیقت، مقدار  $MF_j$  نشان‌گر اهمیت نسبی یک توصیف‌کننده در مقایسه با سایر توصیف‌کننده‌ها در مدل می‌باشد. در مجموعه‌ی توصیف‌کننده‌های انتخاب‌شده بهینه با روش‌های انتخاب متغیر مختلف، توصیف‌کننده‌ی LPRS دارای بیشترین مقدار  $MF_j$  بیشترین تأثیر را نسبت به سایر متغیرها در مدل ارائه شده خواهد داشت. علامت آن، نشان‌گر جهت تغییر در مقادیر عددی متغیر وابسته به عنوان نتیجه‌ای از کاهش یا افزایش در مقادیر توصیف‌کننده است. مقادیر عددی محاسبه‌شده‌ی اثر متوسط ( $MF_j$ ) در جدول ۴ آورده شده‌اند. بر اساس این جدول، مقدار  $MF_j$  نسبت داده شده به توصیف‌کننده‌ی به مراتب بیشتر از سایر توصیف‌کننده‌ها بوده بنابراین انتظار می‌رود این توصیف‌کننده دارای بیشترین تأثیر باشد. بنابراین، معادله‌ی به‌دست آمده برای داده‌های اندیس بازداری این سری از ترکیبات به صورت معادله‌ی ریاضی ۵ قابل بیان می‌باشد:

$$\mathbf{RI} = 406.340(\pm 24.9) + 21.199(\pm 0.6) \mathbf{LPRS} + 305.364(\pm 58.2) \mathbf{MATS1e} - 5515.943(\pm 816.7) \mathbf{Mor29m} - 644.946(\pm 156.6) \mathbf{R7m} + \quad (5)$$

اعداد داخل پرانتز، بیان‌گر خطای استاندارد (انحراف استاندارد) ضرایب و عرض از مبدأ واردشده در مدل خطی بر اساس رگرسیون خطی چندگانه هستند. طبق جدول ۴، هر کدام از توصیف‌کننده‌های وارد شده در مدل، یک خصوصیت ساختاری از مولکول را بیان می‌کنند که نشان‌دهنده‌ی تأثیر بسزای ساختار مولکول‌های مورد بررسی بر اندیس کوتاس آن‌ها بوده و از طرفی ضریب هم‌بستگی میان توصیف‌کننده‌های وارد شده در مدل منتخب، نشان‌دهنده‌ی عدم هم‌بستگی قابل توجه میان این توصیف‌کننده‌ها می‌باشد.

علامت مثبت نسبت داده شده به هر توصیف‌کننده که در جدول ۴ آورده شده است، دلالت بر تقویت مدل به وسیله‌ی آن‌ها داشته در حالی که علامت منفی بیان‌گر رابطه‌ی معکوس بین متغیر وابسته (اندیس کوتاس) و متغیرهای مستقل و توصیف‌کننده‌های مولکولی دارد. بر این اساس، توصیف‌کننده‌های LPRS و MATS1e به ترتیب از گروه‌های توپولوژیکال و خود ارتباطی دوبرعی،

دارای تأثیر سازنده و افزایشی بر مدل خطی ارائه شده هستند. یعنی، با افزایش مقادیر این دو توصیف کننده مقدار اندیس کواتس به عنوان یک متغیر وابسته نیز افزایش می‌یابد. در عین حال، توصیف گرهای Mor29m و R7m+ دارای علامت منفی هستند.

جدول ۴. مشخصات مدل منتخب سری داده‌ی اسانس گیاه میخک زیتنی به روش SW-MLR

No.	Symbol	Descriptor description	Descriptor group	Coefficient	Mean effect	VIF
1	Intercept		-	405.8	-	-
2	LPRS	log of product of row sums (PRS)	Topological	21.2	1.1413	1.6181
3	MATS1e	Moran autocorrelation of lag 1 weighted by Sanderson electronegativity	2D autocorrelations	302.6	-0.0406	1.4854
4	Mor29m	signal 29 / weighted by mass	3D-MoRSE	-5445.8	-0.1065	1.8875
5	R7m+	R maximal autocorrelation of lag 7 / weighted by mass	GETAWAY	-641.7	0.0058	1.4133

این توصیف گرها به ترتیب به خانواده‌های 3D-MoRSE و GETAWAY تعلق دارند. نکته‌ی قابل تأمل دیگر، آن است که قدر مطلق مقدار اثر متوسط ( $MF_j$ ) (رابطه‌ی ۴)، توصیف کننده‌ی LPRS در مقایسه با سایر توصیف کننده‌ها بیشتر است. این بدان معنی است که این توصیف کننده در قیاس با سایر توصیف کننده‌ها دارای بیشترین تأثیر در مدل خطی ارائه شده است. مورد دیگر قابل توجه آن است که مقادیر عددی فاکتور تورم تغییر (VIF)، براساس رابطه‌ی ۳، در گستره‌ی ۱/۴ تا ۱/۹ قرار داشته و این به مفهوم قابل قبول بودن و پایدار بودن مدل ارائه شده می‌باشد. همچنین، جدول ۵ ماتریس هم‌بستگی توصیف کننده‌های وارد شده در مدل نهایی برای سری داده‌ها را نمایش می‌دهد. مقادیر عددی مندرج در این جدول، بیان‌گر ضرایب هم‌بستگی بین هر جفت توصیف کننده‌ی انتخاب شده در مدل خطی نهایی هستند. بر اساس جدول ۵، بیشترین مقدار ضریب هم‌بستگی در این ماتریس برابر با 0.541 و بین توصیف کننده‌های MATS1e و Mor29m می‌باشد. بازه‌ی عددی بسیار کم توصیف کننده‌های انتخابی دلالت بر اعتبار مدل ارائه شده و شاهده‌ی بر رفتار مستقل این متغیرها می‌باشد.

جدول ۵. ماتریس هم‌بستگی ۴ توصیف کننده‌ی وارد شده در مدل منتخب سری داده‌ی اسانس گیاه میخک زیتنی به روش SW-MLR

	LPRS	MATS1e	Mor29m	R7m+
LPRS	1	-0.066	0.541	0.052
MATS1e	-0.066	1	0.346	-0.5
Mor29m	0.541	0.346	1	-0.316
R7m+	0.052	-0.5	-0.316	1

در ادامه هنگامی که از عدم وجود هم‌بستگی بین توصیف‌کننده‌ها اطمینان حاصل شد، با استفاده از معادله‌ی ریاضی ۵ مربوط به مدل ارائه شده، ضرایب رگرسیون موجود در جدول ۴ و نیز مقادیر وارد شده‌ی توصیف‌کننده‌ها در مدل، اندیس‌های بازداری ترکیبات را توسط نرم‌افزار قدرتمند MATLAB محاسبه و پیش‌بینی می‌کنیم. نتایج حاصل شده در این محاسبه در جدول ۶ نشان داده شده است.

جدول ۶. مقادیر تجربی و پیش‌بینی‌شده‌ی اندیس‌های بازداری (RI)، اسانس گیاه میخک زیتی توسط روش رگرسیون خطی چندگانه‌ی مرحله‌ای

No.	Compound name	K.I. (Lit.) <sup>a</sup>	K.I. (Cal.) <sup>b</sup>	Dif. <sup>c</sup>	R.E. <sup>d</sup> %
Training set					
1	$\alpha$ -Thujene	924	944.5	20.5	2.2
2	$\alpha$ -Pinene	930	971.8	41.8	4.5
3	Camphene	938	931.1	-6.9	-0.7
4	Sabinene	958	934.4	-23.6	-2.5
5	Dehydro-1,8-Cineole	973	1046.4	73.4	7.5
6	<i>n</i> -Octanal	973	924.4	-48.6	-5
7	Myrcene	975	996.7	21.7	2.2
8	$\alpha$ -Phellandrene	995	997.9	2.9	0.3
9	$\alpha$ -Terpinene	1002	1002.2	0.2	0
10	<i>p</i> -Cymene	1003	1023.8	20.8	2.1
11	Limonene	1009	998.8	-10.2	-1
12	<i>cis</i> - $\beta$ -Ocimene	1017	979.2	-37.8	-3.7
13	<i>trans</i> - $\beta$ -Ocimene	1027	1029.4	2.4	0.2
14	$\gamma$ -Terpinene	1035	1020.8	-14.2	-1.4
15	<i>trans</i> -Sabinene hydrate	1037	1122.6	85.6	8.3
16	Terpinolene	1064	1003.1	-60.9	-5.7
17	<i>cis</i> -Sabinene hydrate	1066	1130.3	64.3	6
18	<i>n</i> -Nonanal	1073	1067.4	-5.6	-0.5
19	Linalool	1074	1097.3	23.3	2.2
20	<i>allo</i> -Ocimene	1110	1056.3	-53.7	-4.8
21	<i>cis-p</i> -Menthen-1-ol	1110	1128.7	18.7	1.7
22	Hexyl isobutyrate	1127	1213.8	86.8	7.7

23	Terpinene-4-ol	1148	1150.7	2.7	0.2
24	$\alpha$ -Terpineol	1159	1130	-29	-2.5
25	Nerol	1206	1194.2	-11.8	-1
26	Citronellol	1208	1172.7	-35.3	-2.9
27	Hexyl 2-Methylbutyrate	1222	1272.9	50.9	4.2
28	Geraniol	1236	1194.2	-41.8	-3.4
29	$\delta$ -Elemene	1332	1381.4	49.4	3.7
30	$\alpha$ -Cubebene	1345	1385.5	40.5	3
31	$\beta$ -Elemene	1388	1419.6	31.6	2.3
32	$\alpha$ -Gurjunene	1400	1401.5	1.5	0.1
33	$\beta$ -Caryophyllene	1414	1415.9	1.9	0.1
34	Aromadendrene	1428	1405.2	-22.8	-1.6
35	<i>trans</i> - $\alpha$ -Bergamotene	1434	1410.7	-23.3	-1.6
36	$\alpha$ -Humulene	1447	1485.5	38.5	2.7
37	<i>allo</i> -Aromadendrene	1454	1355.4	-98.6	-6.8
38	Germacrene-D	1474	1487	13	0.9
39	Valencene	1477	1446.3	-30.7	-2.1
40	Bicyclogermacrene	1487	1376.9	-110.1	-7.4
41	$\alpha$ - <i>trans-trans</i> -Farnesene	1500	1624.6	124.6	8.3
42	Calamenene	1505	1517.9	12.9	0.9
43	$\delta$ -Cadinene	1505	1458	-47	-3.1
44	Globulol	1566	1570	4	0.3
45	Viridiflorol	1569	1514.4	-54.6	-3.5
46	Guaiol	1575	1567.1	-7.9	-0.5
47	$\delta$ -Cadinol	1618	1668.3	50.3	3.1
48	$\alpha$ -Eudesmol	1634	1599.7	-34.3	-2.1
49	Benzyl benzoate	1701	1660.1	-40.9	-2.4
50	Benzyl salicylate	1790	1748.3	-41.7	-2.3
51	Phenylethyl salicylate	1966	1973.2	7.2	0.4

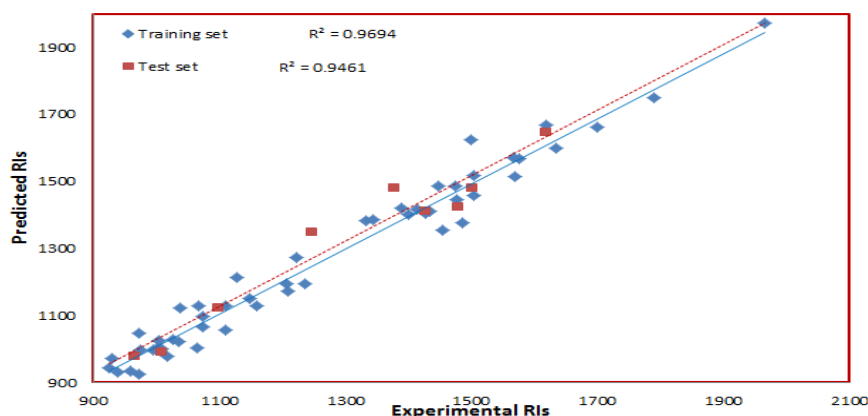
Test set					
1	$\beta$ -Pinene	963	981.2	18.2	1.9
2	$\beta$ -Phellandrene	1005	994.8	-10.2	-1
3	<i>trans-p</i> -Menth-2-en-1-ol	1095	1124.2	29.2	2.7
4	Methyl citronellate	1245	1352.1	107.1	8.6
5	$\alpha$ -Copaene	1375	1483.2	108.2	7.9
6	$\beta$ -Gurjunene	1426	1413.6	-12.4	-0.9
7	$\beta$ -Selinene	1476	1425.3	-50.7	-3.4
8	$\gamma$ -Cadinene	1500	1481.6	-18.4	-1.2
9	T-Cadinol	1616	1649.7	33.7	2.1

A : Retention index in literature: خطای نسبی: D: Relative Error: تفاوت: C: اندیس کواتس محاسبه شده: B: Retention index calculated: اندیس کواتس در منابع علمی:

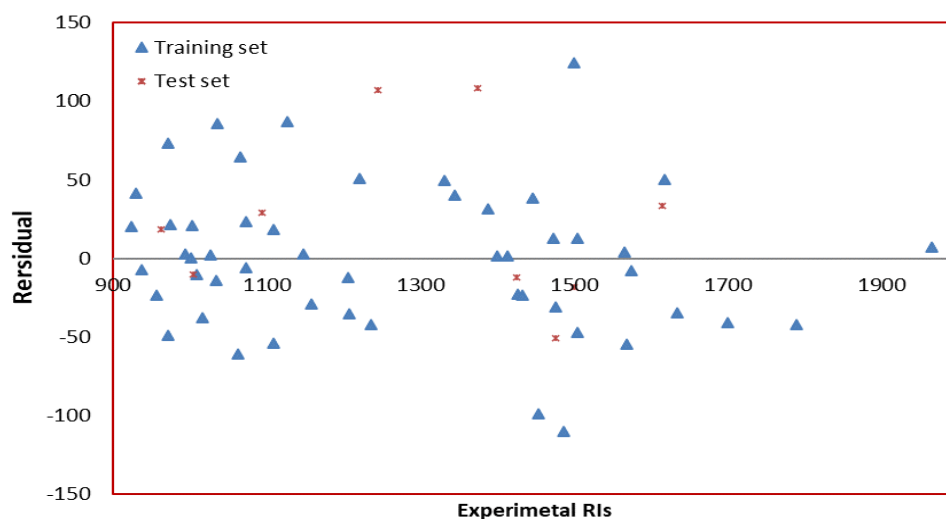
### ۳-۵. نمودارهای پیش‌بینی و باقیمانده

یکی از مواردی که می‌توان با استناد به آن از معتبر بودن مدل انتخابی اطمینان خاطر پیدا کرد، به دست آوردن ضریب اندازه‌گیری  $R^2$  بین مقادیر تجربی یک کمیت و مقادیر نظیر پیش‌گوئی شده برای آن می‌باشد. براساس یک قاعده کلی، هر قدر مقدار ضریب اندازه‌گیری در نمودار پیش‌گوئی به یک نزدیک‌تر باشد، مدل ساخته شده معتبرتر خواهد بود. جهت ارزیابی شمائی اعتبار مدل‌های انتخاب شده، مقادیر پیش‌بینی شده‌ی اندیس‌های بازدارنده را بر حسب مقادیر تجربی مربوطه (شکل ۵) رسم می‌کنیم. براساس شکل ۵، یک هم‌گرایی مطلوب بین مقادیر تجربی و مقادیر پیش‌گوئی شده به چشم می‌خورد. در این شکل، مقدار ضریب اندازه‌گیری برای مجموعه‌های آموزش و آزمون به ترتیب معادل ۰/۹۶۹۴ و ۰/۹۴۶۱ است که دلالت بر اعتبار مدل و توانمندی بالای آن جهت پیش‌گوئی مقادیر اندیس کواتس گسترده‌ی وسیعی از ترکیبات دارد.

مورد دیگری که می‌توان برای اعتبار مدل به آن استناد کرد، رسم مقادیر باقی‌مانده‌ی حاصل از اختلاف مقادیر پیش‌بینی شده‌ی اندیس‌های بازدارنده بر حسب مقادیر تجربی آن می‌باشد. این نمودار، موسوم به نمودار باقی‌مانده در شکل ۶ آورده شده است.



شکل ۵. مقادیر پیش‌بینی شده‌ی اندیس‌های بازدارنده ترکیبات اساس گیاه میخک زینتی بر حسب مقادیر تجربی آن



شکل ۶. مقادیر باقی مانده‌ی حاصل از اختلاف مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازدارنده ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن در نمودار باقی مانده، می‌توان پراکندگی طبیعی نسبتاً یکسان نقاط را حول مقادیر صفر، ناشی از عدم وجود خطای معین در روش و اعتبار بالای مدل دانست.

### ۳-۶. اعتبارسنجی تقاطعی مدل خطی ارائه شده

یکی از رایج‌ترین روش‌های دیگر که پیش‌تر بدان پرداخته شد، اعتبارسنجی تقاطعی (CV) نام دارد. در این روش، در هر مرحله یک یا یک گروه کوچک از داده‌ها کنار گذاشته می‌شود. سپس برای داده‌هایی که باقی مانده، مدلی محاسبه شده و پاسخ از روی مدل محاسبه شده برای یک داده یا گروهی از داده‌ها که کنار گذاشته شده است پیش‌بینی می‌شود. در تکنیک حذف تکی (LOO)، در مراحل متوالی تنها یک ترکیب و مقادیر عددی آن از سری آموزش نهایی حذف شده و مدل‌سازی با بقیه انجام می‌شود. در مرحله بعد، مدل برای ترکیب کنار گذاشته شده، اعمال می‌شود. از طرف دیگر، در تکنیک حذف گروهی (LGO)، در هر مرحله داده‌های مربوط به چند مولکول (پنج ملکول) منحصراً از مجموعه‌ی آموزشی، کنار گذاشته می‌شود. در ادامه، با توجه به داده‌های باقی مانده مدل خطی ساخته می‌شود. سپس، این مدل را برای گروه کنار گذاشته شده اعمال می‌کنیم که مبنای پیش‌بینی اندیس‌های بازدارنده برای آن گروه می‌باشد. در نهایت، با استفاده از یک معادله‌ی ریاضی، تحت مدل ورود اجباری (Enter) و اعمال ضرائب توصیف‌کننده‌های وارد شده در مدل، اندیس‌های بازدارنده ترکیبات توسط نرم‌افزار قدرتمند MATLAB هم در شرایط حذف تک مولکول و در شرایط حذف گروهی مولکول‌ها محاسبه و پیش‌بینی شد. نتایج حاصل از انجام روش اعتبارسنجی تقاطعی با روش‌های حذف تکی (LOO) و حذف گروهی (LGO)، در جدول ۷ نشان داده شده است. طبق قرارداد، در مراجع علمی در مطالعات ارتباط کمی ساختار-فعالیت یا ارتباط کمی ساختار-ویژگی شامل اعتبارسنجی تقاطعی، ضریب اندازه‌گیری یا مجذور ضریب هم‌بستگی بین مقادیر تجربی و مقادیر پیش‌گوئی‌شده، با علامت  $Q^2$  نمایش داده می‌شود. بر اساس مقالات مرجع و پراستناد در پایگاه‌های علمی، مقادیر بالا و قابل توجه  $Q^2$ ، بیانگر استحکام [۱۴]، پایداری، صحت و قدرت پیش‌گوئی بالای مدل خطی ارائه‌شده دارند [۱۵]-

۱۸]. براساس نتایج بدست آمده، مقدار  $Q^2$  برای روش‌های حذف تکی (LOO) و حذف گروهی (LGO)، به ترتیب برابر ۰/۹۵۸ و ۰/۹۵۵ است. این مقادیر، نشان‌گر اعتبار تقاطعی مدل ارائه شده هستند.

### ۳-۷. نتایج آماری مدل ارائه شده اسانس گیاه میخک زینتی به روش SW-MLR

در این بخش، مقادیر تجربی و پیش‌بینی شده‌ی اندیس‌های بازداری را به‌طور جداگانه در یک فایل متنی قرار داده و با استفاده از یک برنامه‌ی جامع در نرم‌افزار MATLAB به‌همراه نتایج حاصل از نرم‌افزار SPSS، پارامترهای آماری شاخص مربوطه را محاسبه می‌کنیم. نتایج حاصل از این محاسبه، در جدول ۷ درج شده است.

جدول ۷. پارامترهای آماری جامع به‌دست آمده برای مجموعه‌های آموزش و آزمون اسانس گیاه *Pittosporum undulatum* به روش SW-MLR

شماره	سری منتخب	پارامتر آماری	مقدار عددی
۱	سری آموزش	$R^2$ (الف)	0.969
	سری آزمون		0.946
۲	سری آموزش	Adjusted $R^2$ (ب)	0.967
	سری آزمون		3.6
۲	سری آموزش	REP (%) (ج)	4.3
	سری آزمون		45.4
۳	سری آموزش	RMSEP (د)	56.4
	سری آزمون		45.8
۴	سری آموزش	SEP (ه)	59.9
	سری آزمون		3.5
۵	سری آموزش	RSEP (و)	4.3
	سری آزمون		82.8
۶	سری آموزش	MAE (ز)	218.9
	سری آزمون		362.8
۷	سری آموزش	$F$ (ح)	16.3
	سری آزمون		

الف) ضریب اندازه‌گیری، ب) ضریب اندازه‌گیری تنظیم شده، ج) درصد خطای نسبی پیشگونی، د) جذر میانگین مجذور خطای پیشگونی، ه) خطای استاندارد پیشگونی، و) خطای استاندارد نسبی پیشگونی، ز)

خطای مطلق میانگین، ح) آماره‌ی  $F$

در جدول ۷، مقادیر عددی ضریب اندازه‌گیری تنظیم شده ( $R^2$  Adjusted) و آماره‌ی  $F$  به‌توسط روابط ۶ و ۷ قابل تعیین هستند.

$$R^2_{Adj} = 1 - (1 - R^2) \left( \frac{n-1}{n-p-1} \right) \quad (۶)$$

$$F = \frac{MSR}{MSE} \quad (۷)$$

در رابطه‌ی (۶)،  $R^2$ ،  $n$  و  $p$  به ترتیب بیان‌گر ضریب اندازه‌گیری، تعداد ترکیبات در سری آموزش و تعداد توصیف‌کننده‌ها هستند. در

منابع علمی،  $R^2_{Adj}$  به‌روش مشابهی با  $R^2$  (ضریب اندازه‌گیری) با در نظر گرفتن تعداد درجات آزادی بیان می‌شود. پارامتر  $R^2_{Adj}$

با تقسیم کردن باقیمانده‌ی مجموع مربعات و مجموع مربعات کل بر درجات آزادی مربوطه تنظیم می‌شود. مقدار  $R^2_{Adj}$  در صورتی - که متغیر (توصیف‌کننده‌ی) اضافه‌شده به معادله، واریانس توضیح داده نشده را کاهش ندهد کم خواهد شد [۱۹]. بنابراین،  $R^2_{Adj}$  معیاری مناسب برای مقایسه‌ی مدل‌های دارای تعداد متفاوتی از متغیرهای پیشگوئی‌کننده (توصیف‌کننده‌ها) می‌باشد. در علم آمار، آماره‌ی  $F$  اساساً برای ارزیابی دقت مدل‌های پیشنهادی به‌طور وسیع مورد استفاده قرار می‌گیرد. در رابطه‌ی تعیین آماره‌ی  $F$  ( $V$ )،  $MSR$  و  $MSE$  به ترتیب دلالت بر میانگین مربعات رگرسیون و میانگین مربعات خطا دارند [۲۰]. این دو پارامتر به-توسط روابط ۸ و ۹ قابل تعیین هستند.

$$MSR = \frac{SSR}{m} \quad (۸)$$

$$MSE = \frac{SSE}{n-m-1} \quad (۹)$$

در معادلات ۸ و ۹،  $SSR$  و  $SSE$  نشان‌گر مجموع مربعات رگرسیون و مجموع مربعات خطا هستند.  $m$  و  $n$  هم به ترتیب بیانگر تعداد توصیف‌کننده‌ها (۴) و تعداد ترکیبات در سری آموزش (۵۱) می‌باشند. بعلاوه، پارامتر  $SSE$ ، موسوم به مجموع مربعات خطاها بر اساس رابطه‌ی ۱۰، قابل بیان می‌باشد.

$$SSE = \sum_{i=1}^n (y_{exp} - y_{pred})^2 \quad (۱۰)$$

تأمل روی رابطه‌ی ۱۰ نشان می‌دهد که رابطه‌ی  $SSE$  دقیقاً با رابطه‌ی PRESS برابر می‌باشد. مجموع مربعات خطا، برطبق رابطه‌ی ۱۱، قابل محاسبه است.

$$SSR = \sum_{i=1}^n (y_{pred} - \bar{y}_{pred})^2 \quad (۱۱)$$

اساساً بین  $SSR$  و مجموع مربعات کل ( $SST$ )، رابطه‌ی ۱۲ برقرار است.

$$SST = SSR + SSE \quad (۱۲)$$

رابطه‌ی مجموع مربعات کل با مقادیر تجربی و میانگین مقادیر تجربی، عبارت است از:

$$SST = \sum_{i=1}^n (y_{exp} - \bar{y}_{exp})^2 \quad (۱۳)$$

در روابط ۱۰ تا ۱۳،  $y_{pred}$ ،  $\bar{y}_{pred}$ ،  $y_{exp}$  و  $\bar{y}_{exp}$  به ترتیب دلالت بر مقدار تجربی کمیت، مقدار پیشگوئی‌شده، میانگین مقادیر پیشگوئی‌شده و میانگین مقادیر تجربی دارند.

### ۳-۸. آزمون $Y$ تصادفی یا آزمون هم‌بستگی شانسی

آزمون  $Y$  تصادفی، به عنوان یک معیار مناسب جهت تأیید عدم وجود هم‌بستگی شانسی بین متغیرهای واردشده در مدل و استحکام مدل در مدل‌سازی‌های  $QSAR$ ،  $QSPR$ ،  $QSRR$  و.... در نظر گرفته می‌شود [۲۱]. در این آزمون، بردار متغیر وابسته (اندیس کواتس) به طور تصادفی در هم آمیخته‌شده یا در هم ریخته‌شده و مدل‌های جدید رابطه‌ی کمی ساختار-فعالیت یا رابطه‌ی کمی



ساختار-ویژگی، با استفاده از ماتریس متغیرهای مستقل اولیه (توصیف‌کننده‌های منتخب) ساخته می‌شود. این در هم‌ریختگی، هرگونه رابطه‌ی منطقی بین توصیف‌کننده‌ها و متغیر وابسته را از بین می‌برد. بنابراین، انتظار می‌رود که با استفاده از این رویکرد پس از چندین با تکرار، مدل‌های ساخته شده دارای مقادیر کم و قابل اغماض  $R^2$  و  $Q^2$  باشند. در جدول ۸، برخی از پارامترهای آماری حاصل از انجام آزمون  $Y$  تصادفی پس از ۱۰ بار تکرار آورده شده‌اند.

جدول ۸. پارامترهای آماری شاخص حاصل از انجام آزمون  $Y$  تصادفی

PRESS <sup>(الف)</sup>	S <sub>PRESS</sub> <sup>(ب)</sup>	SST <sup>(ج)</sup>	R <sup>2</sup> CV <sup>(د)</sup>	PRESS/SST <sup>(ه)</sup>	REP <sup>(و)</sup>	RMSEP <sup>(ز)</sup>
3447433.0582	273.7594	1240222.9990	0.0894	2.7797	272132.0672	259.9937
4051128.4856	296.7625	512329.1367	0.0015	7.9073	319786.3311	281.8402
4300086.3960	305.7452	470887.4879	0.0238	9.1319	339438.4692	290.3712
3164133.5830	262.2699	624770.4901	0.0935	5.0645	249769.0886	249.0820
4198907.2720	302.1268	302585.7967	0.0507	13.8767	331451.6327	286.9347
3969905.2911	293.7725	1089984.7350	0.0215	3.6422	313374.7676	279.0005
4573501.1133	315.3156	624436.3183	0.0294	7.3242	361021.1689	299.4603
4171664.6610	301.1451	829740.2613	0.0009	5.0277	329301.1666	286.0024
4076430.0950	297.6878	395942.7187	0.0111	10.2955	321783.5792	282.7189
4667150.0829	318.5275	352433.5929	0.1546	13.2426	368413.5931	302.5107

(الف) مجموع مربعات باقی‌مانده‌ی پیش‌بینی (ب) عدم قطعیت مجموع مربعات باقی‌مانده‌ی پیش‌بینی (ج) مجموع مربعات کل (د) مجذور ضریب هم‌بستگی اعتبار تقاطعی (ه) مجموع مربعات باقی‌مانده‌ی پیش‌بینی تقسیم بر مجموع مربعات کل (و) خطای نسبی پیش‌بینی (ز) جذرمیانگین مجذور خطای پیش‌بینی

### ۳-۹. نتایج محاسبات و مدل‌های خطی حاصل از سری داده میخک زینتی به روش GA-MLR

پس از انجام محاسبات با استفاده از الگوریتم توانمند GA-MLR با تغییر متوالی در پارامترهای مربوطه در دستورهای متنی در محیط نرم‌افزار متلب، مجموعه‌های شامل ۳ تا ۷ توصیف‌کننده‌ی مولکولی حاصل شد. در جداول متوالی ۹ تا ۱۶، مشخصات و جزئیات کامل این مدل‌ها آورده شده‌اند.

جدول ۹. فهرست توصیف‌کننده‌های منتخب در مدل‌های خطی سری داده‌ی ترکیبات اسانس میخک زینتی به روش GA-MLR

Model number	Descriptor Name(s)	$R^2$		$Q^2$	
		Training set	Test set	$Q^2_{LOO}$	$Q^2_{LGO}$
1	-	-	-	-	-
2	-	-	-	-	-
3	LPRS; Mor18v; R3v	0.9338	0.9775	0.921	0.912
4	LPRS; MATS2e; GATS8e; Mor24e	0.9397	0.9761	0.901	0.882
5	LPRS; MATS2e; GATS8e; Mor11e; R4p+	0.9513	0.9678	0.922	0.910
6	LPRS; MATS2e; SP12; RDF060m; Mor32u; Mor22m	0.946	0.9865	0.920	0.898
7	LPRS; MATS4p; Mor18v; Mor12e; R7m+; R3v; R1e+	0.962	0.9739	0.920	0.912

جدول ۱۰. مشخصات مدل‌های منتخب سری داده‌ی ترکیبات اساس میخک زیتنی به روش GA-MLR

Linear Model 1 (3 molecular descriptors)						
No.	Symbol	Descriptor description	group Descriptor	Coefficient	Mean effect	VIF
1	Intercept	-	-	388.7963	-	-
2	<b>LPRS</b>	log of product row sums (PRS)	topological	17.6513	<b>0.934</b>	1.446
3	Mor18v	3D-MoRSE - signal 18 / weighted by atomic van der Waals volumes	3D-MoRSE	- 143.4047	0.036	1.249
4	R3v	R autocorrelation of lag 3 / weighted by atomic van der Waals volumes	GETAWAY	47.0196	0.031	1.475
Linear Model 2 (4 molecular descriptors)						
1	Intercept	-	-	430.6207	-	-
2	<b>LPRS</b>	log of product row sums (PRS)	topological	19.0013	<b>1.056</b>	1.058
3	MATS2e	Moran autocorrelation - lag 2 / weighted by atomic Sanderson electronegativities	2D autocorrelations	-386.9065	-0.037	1.088
4	GATS8e	Geary autocorrelation - lag 8 / weighted by atomic Sanderson electronegativities	2D autocorrelations	- 20.9688	-0.018	1.540
5	Mor24e	3D-MoRSE - signal 24 / weighted by atomic Sanderson electronegativities	3D-MoRSE	- 14.0612	<b>0.000</b>	1.550
Linear Model 3 (5 molecular descriptors)						
1	Intercept	-	-	473.1039	-	-
2	<b>LPRS</b>	log of product row sums (PRS)	topological	19.2497	<b>1.127</b>	1.076
3	MATS2e	Moran autocorrelation - lag 2 / weighted by atomic Sanderson electronegativities	2D autocorrelations	- 394.9530	-0.040	1.082
4	GATS8e	Geary autocorrelation - lag 8 / weighted by atomic Sanderson electronegativities	2D autocorrelations	-25.1673	-0.023	1.305
5	Mor11e	3D-MoRSE - signal 11 / weighted by atomic Sanderson electronegativities	3D-MoRSE	- 48.3640	-0.059	1.024
6	R4p+	R maximal autocorrelation of lag 4 / weighted by atomic polarizabilities	GETAWAY	- 34.2877	-0.004	1.371
Linear Model 4 (6 molecular descriptors)						
1	Intercept	-	-	395.0259	-	-
2	<b>LPRS</b>	log of product row sums (PRS)	topological	22.0531	<b>1.175</b>	5.238
3	MATS2e	Moran autocorrelation - lag 2 / weighted by atomic Sanderson electronegativities	2D autocorrelations	- 434.8939	-0.040	1.461
4	SP12	shape profile no. 12	Randic molecular profiles	- 14.7261	-0.057	2.372
5	RDF060m	Radial Distribution Function - 6.0 / weighted by atomic masses	RDF	- 24.0879	-0.034	3.211
6	Mor32u	3D-MoRSE - signal 32 / unweighted	3D-MoRSE	4.7085	-0.001	1.846
7	Mor22m	3D-MoRSE - signal 22 / weighted by atomic masses	3D-MoRSE	- 274.0886	-0.043	1.500
Linear Model 5 (7 molecular descriptors)						
1	Intercept	-	-	622.846	-	-
2	<b>LPRS</b>	log of product row sums (PRS)	topological	17.901	<b>1.291</b>	2.580
3	MATS4p	Moran autocorrelation - lag 4 / weighted by atomic polarizabilities	2D autocorrelations	- 81.381	0.014	1.469
4	Mor18v	3D-MoRSE - signal 18 / weighted by atomic van der Waals volumes	3D-MoRSE	- 80.698	0.028	1.918
5	Mor12e	3D-MoRSE - signal 12 / weighted by atomic Sanderson electronegativities	3D-MoRSE	19.045	-0.036	2.195
6	R7m+	R maximal autocorrelation of lag 7 / weighted by atomic masses	GETAWAY	-3248.092	-0.010	1.315
7	R3v	R autocorrelation of lag 3 / weighted by atomic van der Waals volumes	GETAWAY	- 11.856	-0.011	2.031
8	R1e+	R maximal autocorrelation of lag 1 / weighted by atomic Sanderson electronegativities	GETAWAY	- 548.037	-0.277	1.454

جدول ۱۱. فهرست مدل‌های خطی نهائی حاصل در سری داده‌ی ترکیبات اسانس میخک زیتتی به روش GA-MLR

Model number	Linear Model Form
1	$RI = 388.7963(\pm 47.7) + 17.6513(\pm 0.9) \times LPRS - 143.4047(\pm 43.6) \times Mor18v + 47.0196(\pm 87.0) \times R3v$
2	$RI = 430.6207(\pm 34.4) + 19.0013(\pm 0.7) \times LPRS - 386.9065(\pm 123.7) \times MATS2e - 20.9688(\pm 10.4) \times GATS8e - 14.0612(\pm 40.7) \times Mor24e$
3	$RI = 473.1039(\pm 43.0) + 19.2497(\pm 0.7) \times LPRS - 394.9530(\pm 112.1) \times MATS2e - 25.1673(\pm 8.7) \times GATS8e - 48.3640(\pm 14.8) \times Mor11e - 34.2877(\pm 325.2) \times R4p+$
4	$RI = 395.0259(\pm 43.8) + 22.0531(\pm 1.5) \times LPRS - 434.8939(\pm 138.6) \times MATS2e - 14.7261(\pm 6.6) \times SP12 - 24.0879(\pm 14.1) \times RDF060m + 4.7085(\pm 60.5) \times Mor32u - 274.0886(\pm 91.4) \times Mor22m$
5	$RI = 622.846(\pm 94.6) + 17.901(\pm 0.9) \times LPRS - 81.381(\pm 111.2) \times MATS4p - 80.698(\pm 42.8) \times Mor18v + 19.045(\pm 22.6) \times Mor12e - 3248.092(\pm 637.8) \times R7m+ - 11.856(\pm 80.9) \times R3v - 548.037(\pm 213.3) R1e+$

جدول ۱۲. پارامترهای آماری جامع به دست آمده برای مجموعه‌های آموزش و آزمون در سری داده‌ی ترکیبات اسانس میخک زیتتی به روش GA-MLR در مدل شامل ۳ توصیف‌کننده

شماره	سری منتخب	پارامتر آماری	مقدار عددی
۱	سری آموزش	$R^2$	0.9338
	سری آزمون		0.9775
۲	سری آموزش	Adjusted $R^2$	0.934
	سری آزمون		5.3
۳	سری آموزش	REP (%)	2.7
	سری آزمون		66.7
۴	سری آموزش	RMSEP	35.6
	سری آزمون		67.4
۵	سری آموزش	SEP	37.7
	سری آزمون		5.2
۶	سری آموزش	RSEP	2.7
	سری آزمون		99.9
۷	سری آموزش	MAE	178.9
	سری آزمون		220.9
۸	سری آموزش	$F$	59.9
	سری آزمون		

جدول ۱۳. پارامترهای آماری جامع به دست آمده برای مجموعه‌های آموزش و آزمون در سری داده‌ی ترکیبات اسانس میخک زینتی به روش GA-MLR در

مدل شامل ۴ توصیف‌کننده

شماره	سری منتخب	پارامتر آماری	مقدار عددی
۱	سری آموزش	$R^2$	0.9338
	سری آزمون		0.9775
۲	سری آموزش	Adjusted $R^2$	0.934
	سری آزمون		5.0
۲	سری آموزش	REP (%)	2.7
	سری آزمون		63.7
۳	سری آموزش	RMSEP	35.3
	سری آزمون		64.4
۴	سری آموزش	SEP	37.5
	سری آزمون		4.9
۵	سری آموزش	RSEP	2.7
	سری آزمون		100.8
۶	سری آموزش	MAE	180.6
	سری آزمون		179.1
۷	سری آموزش	$F$	34.9
	سری آزمون		

جدول ۱۴. پارامترهای آماری جامع به دست آمده برای مجموعه‌های آموزش و آزمون در سری داده‌ی ترکیبات اسانس میخک زینتی به روش GA-MLR در

مدل شامل ۵ توصیف‌کننده

شماره	سری منتخب	پارامتر آماری	مقدار عددی
۱	سری آموزش	$R^2$	0.9397
	سری آزمون		0.9761
۲	سری آموزش	Adjusted $R^2$	0.946
	سری آزمون		4.5
۲	سری آموزش	REP (%)	3.2
	سری آزمون		57.3
۳	سری آموزش	RMSEP	41.9
	سری آزمون		57.8
۴	سری آموزش	SEP	44.5
	سری آزمون		4.4
۵	سری آموزش	RSEP	3.2
	سری آزمون		96.5
۶	سری آموزش	MAE	207.7
	سری آزمون		175.7
۷	سری آموزش	$F$	14.4
	سری آزمون		

جدول ۱۵. پارامترهای آماری جامع به دست آمده برای مجموعه‌های آموزش و آزمون در سری داده‌ی ترکیبات اسانس میخک زینتی به روش GA-MLR در

مدل شامل ۶ توصیف‌کننده

شماره	سری منتخب	پارامتر آماری	مقدار عددی
۱	سری آموزش	$R^2$	0.9513
	سری آزمون		0.9678
۲	سری آموزش	Adjusted $R^2$	0.939
	سری آزمون		4.8
۲	سری آموزش	REP (%)	2.4
	سری آزمون		60.3
۳	سری آموزش	RMSEP	31.0
	سری آزمون		60.9
۴	سری آموزش	SEP	32.9
	سری آزمون		4.7
۵	سری آموزش	RSEP	2.4
	سری آزمون		98.7
۶	سری آموزش	MAE	157.6
	سری آزمون		128.6
۷	سری آموزش	$F$	16.8
	سری آزمون		

جدول ۱۶. پارامترهای آماری جامع به دست آمده برای مجموعه‌های آموزش و آزمون در سری داده‌ی ترکیبات اسانس میخک زینتی به روش GA-MLR در

مدل شامل ۷ توصیف‌کننده

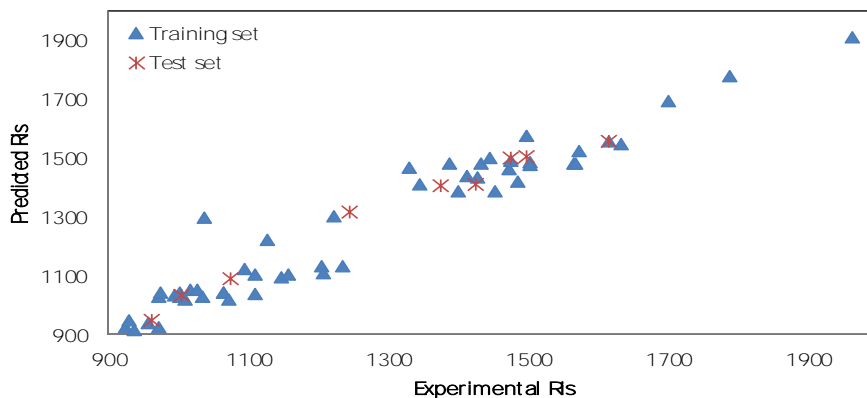
شماره	سری منتخب	پارامتر آماری	مقدار عددی
۱	سری آموزش	$R^2$	0.962
	سری آزمون		0.9739
۲	سری آموزش	Adjusted $R^2$	0.956
	سری آزمون		4.0
۲	سری آموزش	REP (%)	2.9
	سری آزمون		50.6
۳	سری آموزش	RMSEP	37.3
	سری آزمون		51.1
۴	سری آموزش	SEP	39.5
	سری آزمون		3.9
۵	سری آموزش	RSEP	2.8
	سری آزمون		87.2
۶	سری آموزش	MAE	188.2
	سری آزمون		155.6
۷	سری آموزش	$F$	4.8
	سری آزمون		

جهت درک عمیق‌تر و مقایسه‌ی کامل‌تر مدل‌های ساخته‌شده، منحنی‌های پیش‌گوئی و باقی‌مانده برای مدل‌های خطی شامل ۳، ۴، ۵،

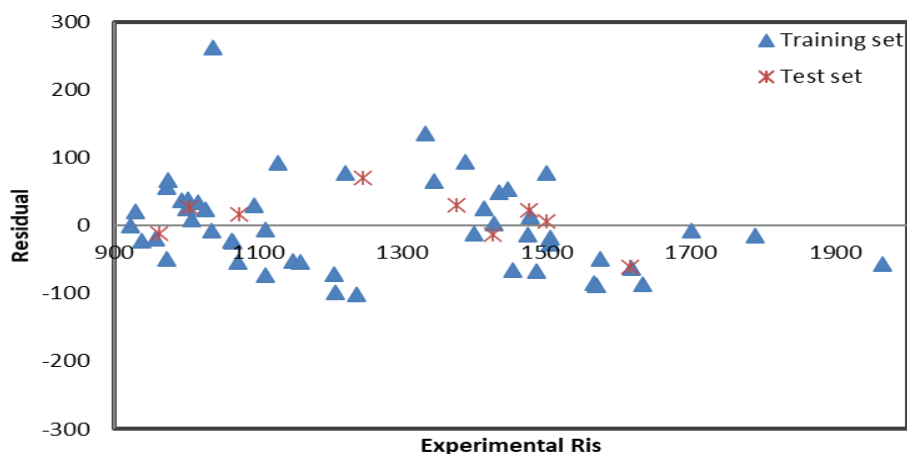
۶ و ۷ توصیف‌کننده‌ی مولکولی به ترتیب در شکل‌های ۱۰، ۱۱، ۱۲، ۱۳، ۱۴، ۱۵، ۱۶، ۱۷، ۱۸ و ۱۹ آورده شده است. نهایتاً، در

جداول ۱۸، ۱۹، ۲۰، ۲۱ و ۲۲، ماتریس هم‌بستگی ۳، ۴، ۵، ۶ و ۷ توصیف‌کننده‌ی واردشده در مدل‌های خطی به روش GA-MLR

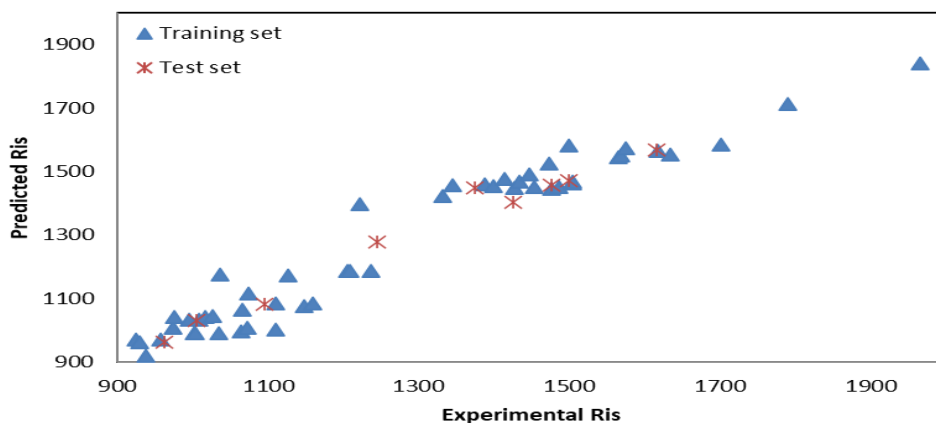
نمایش داده شده است. براساس این جداول، مقادیر ناچیز ضرائب هم‌بستگی بین هر جفت توصیف کننده دلالت بر رفتار مستقل آن‌ها و مناسب بودن مدل‌های ارائه شده دارد.



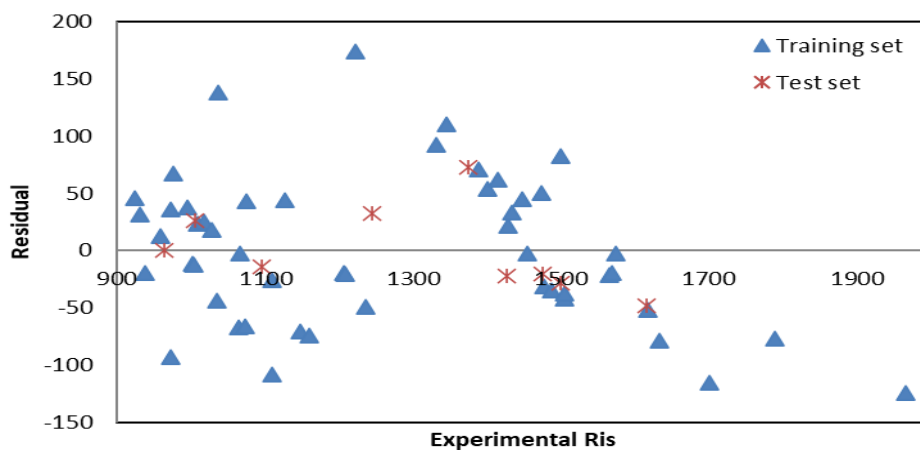
شکل ۱۰. مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازدارنده ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن به‌روش GA-MLR در مدل شامل ۳ توصیف کننده



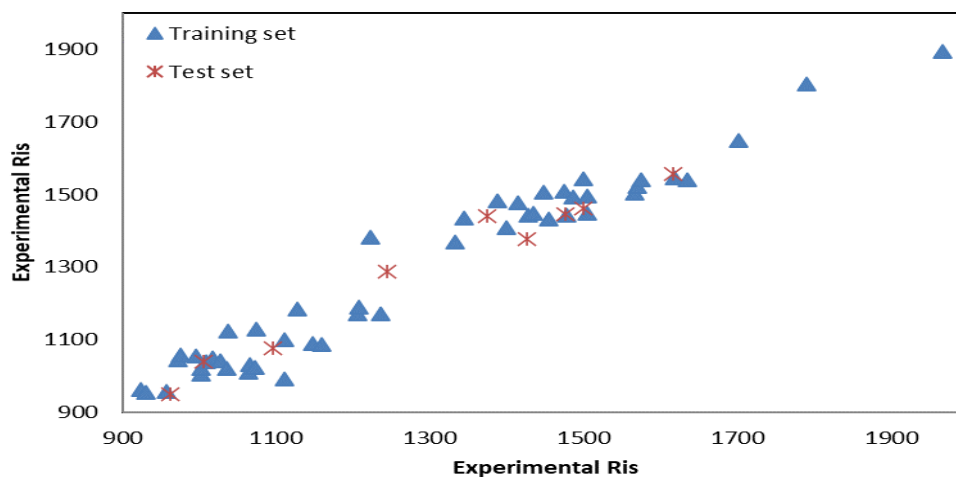
شکل ۱۱. مقادیر باقی‌مانده‌ی حاصل از اختلاف مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازدارنده ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن به‌روش GA-MLR در مدل شامل ۳ توصیف کننده



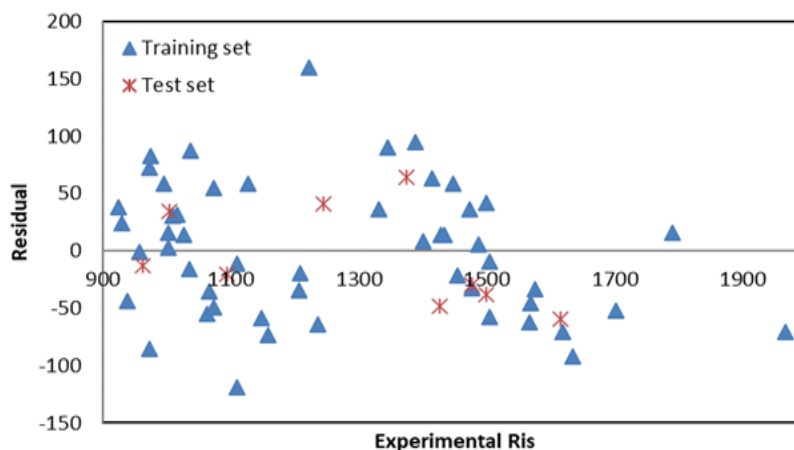
شکل ۱۲. مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازدارنده ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن به‌روش GA-MLR در مدل شامل ۴ توصیف کننده



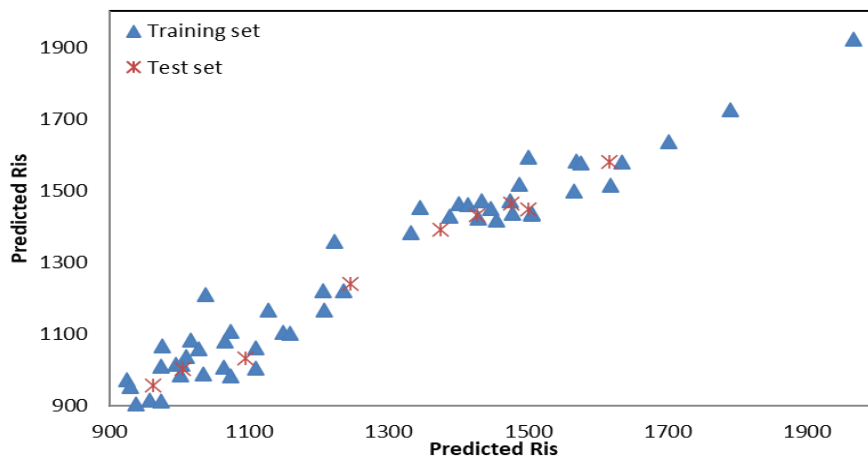
شکل ۱۳. مقادیر باقی مانده‌ی حاصل از اختلاف مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازداری ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن به روش GA-MLR در مدل شامل ۴ توصیف‌کننده



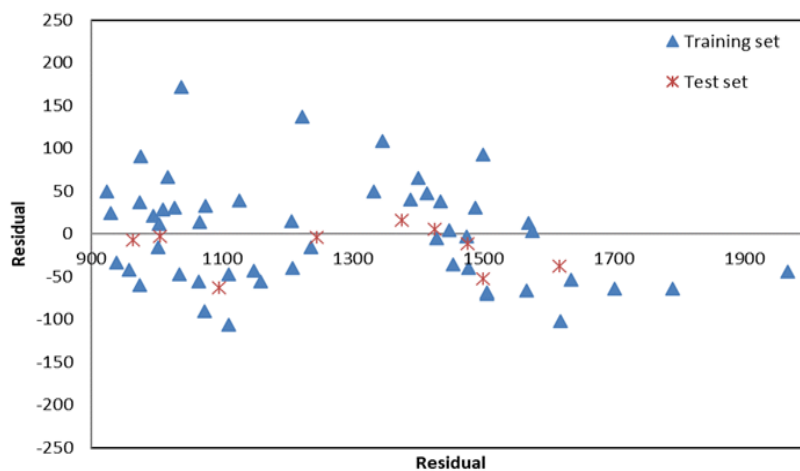
شکل ۱۴. مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازداری ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن به روش GA-MLR در مدل شامل ۵ توصیف‌کننده



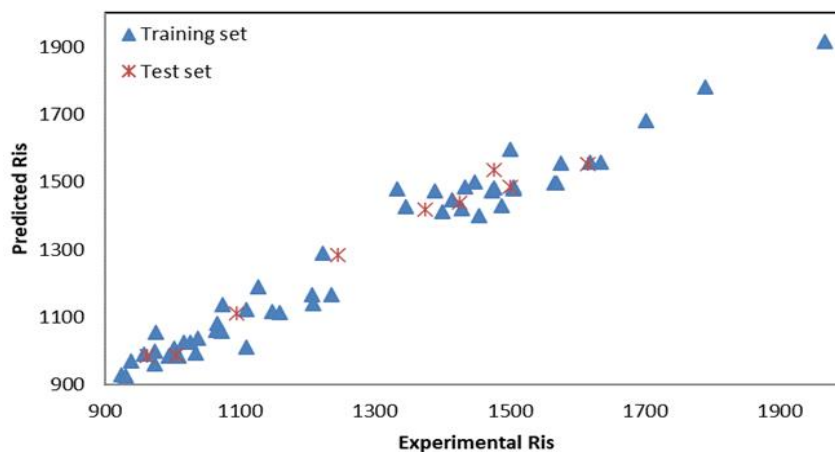
شکل ۱۵. مقادیر باقی مانده‌ی حاصل از اختلاف مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازداری ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن به روش GA-MLR در مدل شامل ۵ توصیف‌کننده



شکل ۱۶. مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازدارنده‌ی ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن به‌روش GA-MLR در مدل شامل ۶ توصیف‌کننده

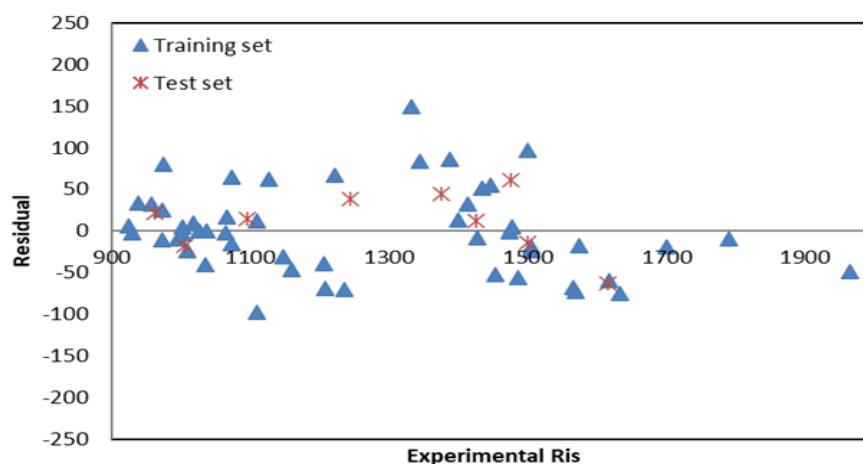


شکل ۱۷. مقادیر باقی‌مانده‌ی حاصل از اختلاف مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازدارنده‌ی ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن به‌روش GA-MLR در مدل شامل ۶ توصیف‌کننده



شکل ۱۸. مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازدارنده‌ی ترکیبات اسانس گیاه میخک زینتی برحسب مقادیر تجربی آن به‌روش GA-MLR در مدل شامل ۷ توصیف‌کننده





شکل ۱۹. مقادیر باقی مانده‌ی حاصل از اختلاف مقادیر پیش‌بینی‌شده‌ی اندیس‌های بازدارنده‌ی ترکیبات اسانس گیاه میخک زینتی بر حسب مقادیر تجربی آن به روش GA-MLR در مدل شامل ۷ توصیف‌کننده

جدول ۱۸. ماتریس هم‌بستگی ۳ توصیف‌کننده‌ی وارد شده در مدل منتخب سری داده‌ی ترکیبات اسانس گیاه میخک زینتی به روش GA-MLR

	LPRS	Mor18v	R3v
LPRS	1	-0.216	0.439
Mor18v	-0.216	1	0.256
R3v	0.439	0.256	1

جدول ۱۹. ماتریس هم‌بستگی ۴ توصیف‌کننده‌ی وارد شده در مدل منتخب سری داده‌ی ترکیبات اسانس گیاه میخک زینتی به روش GA-MLR

	LPRS	MATS2e	GATS8e	Mor24e
LPRS	1	0.216	0.035	-0.014
MATS2e	0.216	1	-0.055	0.103
GATS8e	0.035	-0.055	1	0.577
Mor24e	-0.014	0.103	0.577	1

جدول ۲۰. ماتریس هم‌بستگی ۵ توصیف‌کننده‌ی وارد شده در مدل منتخب سری داده‌ی ترکیبات اسانس گیاه میخک زینتی به روش GA-MLR

	LPRS	MATS2e	GATS8e	Mor11e	R4p+
LPRS	1	0.216	0.035	-0.101	0.124
MATS2e	0.216	1	-0.055	-0.020	0.194
GATS8e	0.035	-0.055	1	0.061	-0.474
Mor11e	-0.101	-0.020	0.061	1	-0.123
R4p+	0.124	0.194	-0.474	-0.123	1

جدول ۲۱. ماتریس هم‌بستگی ۶ توصیف‌کننده‌ی وارد شده در مدل منتخب سری داده‌ی ترکیبات اسانس گیاه میخک زینتی به روش GA-MLR

	LPRS	MATS2e	SP12	RDF060m	Mor32u	Mor22m
LPRS	1	0.216	0.301	0.794	-0.469	0.273
MATS2e	0.126	1	-0.398	0.305	-0.258	0.234
SP12	0.301	-0.398	1	0.184	0.230	-0.365
RDF060m	0.794	0.305	0.184	1	-0.245	0.205
Mor32u	-0.469	-0.258	0.230	-0.245	1	-0.193
Mor22m	0.273	0.234	-0.365	0.205	-0.193	1

جدول ۲۲. ماتریس هم‌بستگی ۷ توصیف‌کننده‌ی وارد شده در مدل منتخب سری داده‌ی ترکیبات اسانس گیاه میخک زینتی به روش GA-MLR

	LPRS	MATS4p	Mor18v	Mor12e	R7m+	R3v	R1e+
LPRS	1	0.236	-0.216	-0.451	0.054	0.439	-0.518
MATS4p	0.263	1	-0.125	-0.142	-0.145	-0.256	-0.226
Mor18v	-0.216	-0.125	1	-0.366	0.238	0.256	0.026
Mor12e	-0.451	-0.142	-0.366	1	0.143	-0.506	0.334
R7m+	0.054	-0.145	0.238	0.143	1	0.029	-0.083
R3v	0.439	-0.256	0.256	-0.506	0.029	1	-0.312
R1e+	-0.518	-0.226	-0.026	0.334	-0.083	-0.312	1

#### ۴. نتیجه‌گیری

در این مقاله، به صورت جامع به مدل‌سازی‌های خطی انجام شده روی ترکیبات طبیعی تشکیل‌دهنده روغن اسانسی استحصال شده از گیاه پرداخته شده است. این مدل‌سازی‌ها، با استفاده از روش‌های SW-MLR و GA-MLR انجام پذیرفته‌اند. جهت نیل به روابط خطی هدفمند و مستحکم پس از انجام مراحل محاسبه و انتخاب توصیف‌کننده‌ها و تجزیه و تحلیل آماری آن‌ها، الگوی پراکنش آنالیز جزء اصلی (PCA) و تقسیم سری داده به مجموعه‌های آموزشی و آزمون به صورت منطقی انجام شد. نتایج حاصل از انجام این نوع آنالیز مشخص نمود که ترکیبات طبیعی با شماره‌های ۳۰، ۵۸، ۵۹ و ۶۰ در مجموعه‌ی داده‌های اصلی با فاصله محسوس نسبت به خوشه‌ی دربرگیرنده‌ی سایر ترکیبات طبیعی شناسایی شده قرار گرفته‌اند. رگرسیون خطی چندگانه مرحله‌ای منجر به مجموعه‌ای شامل ۱۴ مدل رگرسیون در نرم‌افزار SPSS شد. همچنین، ترکیبات شماره‌ی ۱۵ و ۴۵ در سری آموزش دارای مقادیر لورج بیشتر از مقدار هشدار ( $h^*$ ) و می‌توانند جزء ترکیبات انحرافی باشند. خوشبختانه، در دو مورد اخیر مقادیر اطلاعات عددی پیشگونی شده مناسب و قابل قبول هستند. لذا، این دو ترکیب لورج خوب خوانده شده و از حوزه‌ی کاربرد و سری داده‌ی اولیه حذف نمی‌شوند. نکته حائز اهمیت اینکه در کلیه مدل‌های خطی ارائه شده، هم‌گرایی‌های مطلوب بین مقادیر تجربی و مقادیر پیشگونی شده حاصل شدند.

از طرف دیگر، در نمودارهای باقیمانده‌ی ترسیم شده مربوط به تفاضل مقادیر حقیقی و مقادیر پیشگونی شده به توسط مدل‌های خطی مربوطه، پراکنندگی طبیعی نسبتاً یکسان نقاط را حول مقادیر صفر را می‌توان ناشی از عدم وجود خطای معین در روش و اعتبار بالای مدل دانست. بعلاوه، پس از انجام روش اعتبارسنجی تقاطعی (CV)، مقادیر عددی محاسبه شده‌ی ضرایب اندازه‌گیری (یا مجذور ضریب هم‌بستگی بین مقادیر تجربی و مقادیر پیش‌گونی شده) نشانگر پایداری، صحت و قدرت پیشگونی بالای مدل‌های خطی مربوطه هستند.

در فاز دوم این مطالعه‌ی نظری با استفاده از روش انتخاب متغیر توانمند الگوریتم ژنتیک مبتنی بر نظریه‌ی تکامل طبیعی داروین همراه با رگرسیون خطی چندگانه (GA-MLR)، با تغییر متوالی در پارامترهای مربوطه در دستورهای متنی در محیط نرم‌افزار متلب، مجموعه‌های شامل ۳ تا ۷ توصیف‌کننده‌ی مولکولی حاصل شدند که هر یک به تنهایی از استحکام و قدرت پیشگویی قابل ملاحظه

ای جهت پیش‌بینی مقادیر عددی شاخص‌های بازدارنده کوآتس ترکیبات طبیعی برخوردارند. به عنوان پیشنهاداتی تکمیلی، برای انجام مدل‌سازی‌های اینچینی در آینده می‌توان از روش‌های غیرخطی مانند شبکه‌های عصبی مصنوعی و همچنین از الگوریتم‌های دیگر مانند الگوریتم کرم شب‌تاب، الگوریتم مورچه و یا جنگل تصادفی نیز می‌توان بهره‌برداری و توانمندی مدل‌های مختلف را با در نظر گرفتن شاخص‌های آماری مرتبط با هر نوع مدل‌سازی تعیین و با یکدیگر مقایسه نمود.

## ۵. مراجع

- [1] Zargari, A. (1991). Medicinal Plants. Tehran University Publication, Tehran.
- [2] Mozaffarian, V. (1996). A Dictionary of Iranian Plant Names. Farhang Moaser Press, Iran.
- [3] Ferreira, N. J., de Sousa, I. G. M., Luis, T. C., Currais, A. J. M., Figueiredo, A. C., Costa, M. M., Lima, A. S. B., Santos, P. A. G., Barroso, J. G., Pedro, L. G. & Scheffer, J. J. C. (2007). *Pittosporum undulatum* Vent. grown in Portugal: Secretory structures, seasonal variation and enantiomeric composition of its essential oil. *Flavour and Fragrance Journal*, 22, 1-9.
- [4] Wold, S., Esbensen, K. & Geladi, P. (1987). Principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, 2, 37-52.
- [5] Heberger, K. (1999). Evaluation of polarity indicators and stationary phases by principal component analysis in gas-liquid chromatography *Chemometrics and Intelligent Laboratory Systems*, 47, 41-49.
- [6] Heberger, K., Milczewska, K. & Voelkel, A. (2005). Principal component analysis of polymer-solvent and filler-solvent interactions by inverse gas chromatography. *Colloids and Surfaces a-Physicochemical and Engineering Aspects*, 260, 29-37.
- [7] Heberger, K. & Gorgeny, M. (1999). Principal component analysis of Kovats indices for carbonyl compounds in capillary gas chromatography. *Journal of Chromatography A*, 845, 21-31.
- [8] OECD (2007). Guidance Document on the Validation of (Quantitative) Structure-Activity Relationships [(Q)SAR] Models. Organisation for Economic Co-Operation and Development, Paris.
- [9] Netzeva, T. I., Worth, A. P., Aldenberg, T., Benigni, R., Cronin, M. T. D., Gramatica, P., Jaworska, J. S., Kahn, S., Klopman, G., Marchant, C. A., Myatt, G., Nikolova-Jeliazkova, N., Patlewicz, G. Y., Perkins, R., Roberts, D. W., Schultz, T. W., Stanton, D. T., van de Sandt, J. J. M., Tong, W. D., Veith, G. & Yang, C. H. (2005). Current status of methods for defining the applicability domain of (quantitative) structure-activity relationships - The report and recommendations of ECVAM Workshop 52. *Atla-Alternatives to Laboratory Animals*, 33, 155-173.
- [10] Eriksson, L., Jaworska, J., Worth, A. P., Cronin, M. T. D., McDowell, R. M. & Gramatica, P. (2003). Methods for reliability and uncertainty assessment and for applicability evaluations of classification- and regression-based QSARs. *Environmental Health Perspectives*, 111, 1361-1375.
- [11] Jaworska, J., Nikolova-Jeliazkova, N. & Aldenberg, T. (2005). QSAR applicability domain estimation by projection of the training set in descriptor space: A review. *Atla-Alternatives to Laboratory Animals*, 33, 445-459.
- [12] Agrawal, V. K. & Khadikar, P. V. (2001). QSAR prediction of toxicity of nitrobenzenes. *Bioorganic & Medicinal Chemistry*, 9, 3035-3040.
- [13] Pournasheer, E., Riahi, S., Ganjali, M. R. & Norouzi, P. (2010). Quantitative structure-activity relationship (QSAR) study of interleukin-1 receptor associated kinase 4 (IRAK-4) inhibitor activity by the genetic algorithm and multiple linear regression (GA-MLR) method. *Journal of Enzyme Inhibition and Medicinal Chemistry*, 25, 844-853.

- [14] Beheshti, A., Pourbasheer, E., Nekoei, M. & Vahdani, S. (2012). QSAR modeling of antimalarial activity of urea derivatives using genetic algorithm–multiple linear regressions. *Journal of Saudi Chemical Society*, 20, 282-290.
- [15] Gramatica, P. (2007). Principles of QSAR models validation: Internal and external. *QSAR & Combinatorial Science*, 26, 694-701.
- [16] Noorizadeh, H. & Farmany, A. (2012). Quantitative structure-retention relationship for retention behavior of organic pollutants in textile wastewaters and landfill leachate in LC-APCI-MS. *Environmental Science and Pollution Research*, 19, 1252-1259.
- [17] Noorizadeh, H., Noorizadeh, M. & Farmany, A. (2012). Advanced QSRR models of toxicological screening of basic drugs in whole blood by UPLC-TOF-MS. *Medicinal Chemistry Research*, 21, 4357-4368.
- [18] Jalali-Heravi, M. & Asadollahi-Baboli, M. (2009). Quantitative structure-activity relationship study of serotonin (5-HT<sub>7</sub>) receptor inhibitors using modified ant colony algorithm and adaptive neuro-fuzzy interference system (ANFIS). *European Journal of Medicinal Chemistry*, 44, 1463-1470.
- [19] Agirre-Basurko, E., Ibarra-Berastegi, G. & Madariaga, I. (2006). Regression and multilayer perceptron-based models to forecast hourly O<sub>3</sub> and NO<sub>2</sub> levels in the Bilbao area. *Environmental Modelling & Software*, 21, 430-446.
- [20] Larose, D. T. (2006). *Data Mining Methods and Models*. John Wiley & Sons, Inc Publication, Hoboken, New Jersey.
- [21] Tropsha, A., Gramatica, P. & 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, 69-77.

## Steps, calculations and results of studies of theoretical predictions of quantitative structure retention relationship (QSRR) of *Pittosporum undulatum* essential oil

Majid Mohammadhosseini<sup>1\*</sup>, Mehdi Nekoei<sup>1</sup>

<sup>1</sup>Department of Chemistry, Shahrood Branch, Islamic Azad University, Shahrood, Iran

Submitted: 14 April 2023, Revised: 13 July 2023, Accepted: 21 July 2023

### Abstract

In this article, a detailed description of the linear models capable of predicting the inhibition indices of a large group of natural compounds identified in the essential oil of *Pittosporum undulatum*, as one of the medicinal plants, has been discussed. In this regard, the work is based on quantitative structure retention relationship (QSRR), which is of prime importance in scientific resources to establish a logical and meaningful relationship between the Kovats index as a dependent variable and a group of molecular descriptors as independent variables. In this regard, after drawing the structure of the natural compounds using the Hypercam software and optimizing their molecular structures, Dragon software was used to extract the relevant molecular descriptors. In the next step, after removing irrelevant and redundant descriptors, a group of important and effective descriptors were identified and their linear relationship with the Kovats retention index was discussed and investigated using stepwise multiple linear regression method as well as another variable selection method based upon genetic algorithm feature selection approach. The obtained results indicate the high capability of the presented models to predict the Kovats index of a wide group of natural compounds.

**Keywords:** *Quantitative structure retent relationship (QSRR), multiple linear regression, molecular descriptors, genetic algorithm, Kovats retention index, Pittosporum undulatum.*