



روش‌های محاسباتی و مطالعه ارتباط کمی ساختار- خاصیت جهت پیش‌بینی نقاط ذوب ترکیبات نیترواروماتیک حلقوی با استفاده از توصیف‌کننده‌های شیمیایی و مکانیک کوانتومی: ترکیب محاسبات DFT و QSPR

مهدی نکویی^{۱*}، مهدی مهام^۲، بهزاد چهکندی^۱، سید محسن بابازاده^۱
^۱دانشگاه آزاد اسلامی، واحد شاهرود، دانشکده علوم پایه، گروه شیمی، شاهرود، ایران
^۲دانشگاه آزاد اسلامی، واحد علی‌آبادکتول، گروه شیمی، علی‌آبادکتول، ایران

تاریخ ثبت اولیه: ۱۳۹۵/۰۵/۱۲، تاریخ دریافت نسخه اصلاح شده: ۱۳۹۵/۰۷/۱۸، تاریخ پذیرش قطعی: ۱۳۹۵/۰۸/۲۵

چکیده

روش DFT-B3LYP با پایه 6-31G(d) جهت محاسبه چندین توصیف‌کننده مکانیک کوانتومی در ۶۰ ترکیب نیترواروماتیک حلقوی بکار گرفته شد. مناسب‌ترین توصیف‌کننده‌ها جهت مطالعه ارتباط کمی ساختار- خاصیت در پیش‌بینی نقاط ذوب ترکیبات نیترواروماتیک حلقوی با استفاده از توصیف‌کننده‌های مکانیک کوانتومی و شیمیایی و ماشین بردار پشتیبان (SVM) انجام شد. در ابتدا ساختار ترکیبات رسم و گروه مناسبی از توصیف‌کننده‌های شیمیایی و مکانیک کوانتومی محاسبه شدند. سپس از روش انتخاب مرحله‌ای برای بدست آوردن بهترین توصیف‌کننده‌ها که بیشترین ارتباط را با نقاط ذوب ترکیبات مورد نظر داشتند استفاده گردید. در ابتدا مدل خطی رگرسیون خطی چند-گانه (MLR) ساخته شد. سپس برای به دست آوردن مدل بهتر، از SVM استفاده گردید. داده‌های آماری، برتری روش SVM را نسبت به روش خطی MLR نشان می‌دهد.

واژه‌های کلیدی: ارتباط کمی ساختار- خاصیت، نقطه ذوب، رگرسیون خطی چندگانه، ماشین بردار پشتیبان، توصیف‌کننده‌های مکانیک کوانتومی

۱. مقدمه

جستجو برای ترکیبات پرنانرژی جدید با خواص فیزیکی و شیمیایی مطلوب، کارایی و عملکرد بالا و حساسیت کاهش یافته به محرک‌های خارجی از اهمیت بسیار بالایی در صنایع شیمیایی برخوردار است [۱]. در سال‌های اخیر تلاش‌های مداوم و تحقیقات گسترده‌ای در جهت توسعه ترکیبات پرنانرژی جدید به علت کاربردهای زیادشان به عنوان پیشران‌ها در صنایع نظامی و مواد منفجره انجام گرفته است. اما سنتز ترکیبات

*عهده‌دار مکاتبات: مهدی نکویی

نشانی: شاهرود، دانشگاه آزاد اسلامی، گروه شیمی

تلفن: ۰۲۳-۳۲۳۹۴۲۸۹-۰۲۳ پست الکترونیک: E-Mail: m_nekoei1356@yahoo.com

پرانرژی همراه با مخاطرات و مضرات فراوانی است بعلاوه بررسی و اندازه گیری اطلاعات و خواص ترکیبات پرانرژی، گران و وقت گیر و گاهی اوقات حتی غیرممکن است. بنابراین توانایی تخمین خواص مطلوب مواد پرانرژی بطور مستقیم از روی ساختار مولکولیشان بسیار جذاب و با اهمیت است. این امر می تواند در تصمیم گیری برای اینکه آیا تلاش برای سنتز این ترکیبات ارزش ویژه ایی دارد یا خیر کمک شایانی خواهد کرد [۱]. اخیراً روشهای مختلفی برای پیش بینی خواص ترموشیمی از قبیل گرمای تصعید، گرمای تشکیل و غیره برای دسته های مختلف ترکیبات پرانرژی معرفی شده است [۲-۳].

پیش بینی نقطه ذوب (Mp) ترکیبات پرانرژی دارای اهمیت زیادی است زیرا نقطه ذوب از خواص فیزیکوشیمیایی پایه ایی است که در شناسایی و تعیین خلوص ترکیبات وهمچنین در محاسبه برخی دیگر از خواص فیزیکوشیمیایی از قبیل فشار بخار و حلالیت کاربرد دارد. دینامیک مولکولی جهت شبیه سازی انتقال فاز جامد به مایع در مواد پرانرژی به منظور پیش بینی نقاط ذوبشان بکار رفته است اما این کار به علت سد انرژی آزاد جهت تشکیل فاز جامد-مایع، نسبتاً مشکل و وقت گیر است.

روش ارتباط کمی ساختار- خاصیت^۱ (QSPR) یک روش جایگزین برای تخمین نقاط ذوب ترکیبات پرانرژی براساس توصیف کننده های بدست آمده از ساختار مولکولی آنها می باشد. مزیت این روش این است که فقط به فهم و دانش ساختار شیمیایی احتیاج است و به هیچ خاصیت تجربی دیگری وابسته نیست. هدف از مطالعات QSPR پیدا کردن رابطه ایی است که بین رفتار فیزیکوشیمیایی یک مولکول با پارامترهای ساختاری آن وجود دارد [۴-۱۰]. نتایج این مطالعات علاوه بر شفاف سازی نحوه ارتباط بین خواص مولکولها و ویژگی های ساختاری آنها به پژوهشگران در پیش بینی رفتار مولکولهای جدید براساس رفتار مولکولهای مشابه کمک می کند. کار اصلی در مطالعات ارتباط کمی ساختار-ویژگی، برقراری رابطه و مدلی بین خواص مولکولی مشخص و توصیف گرهای مولکولی توسط آمار و یا روش های دیگر است. باید توجه شود که مدل ممکن است کاملاً پیچیده باشد، و بنابراین اغلب انواعی از تقریب ها به کار می رود. تعدادی از روش هایی که به طور وسیعی در مطالعات مدل سازی ارتباط کمی ساختار-ویژگی استفاده می شوند، روش های خطی هستند که مدل های خطی ساده و قابل تفسیری را می سازند، از قبیل روش حداقل مربعات متداول^۲ (OLS)، رگرسیون اجزاء اصلی^۳ (PCR) و حداقل مربعات جزئی^۴ (PLS). برخی دیگر از روش ها مدل های غیرخطی اند، مانند شبکه های عصبی مصنوعی^۵ و ماشین بردار پشتیبان^۶ که بطور گسترده از آنها استفاده می گردد [۱۱-۱۷].

سهم اصلی برای گسترش استفاده از مدل های ارتباط کمی ساختار-ویژگی مربوط است به توسعه توصیف گرهای ساختاری و معادلات ریاضی که خواص فیزیکی و شیمیایی را به ساختار شیمیایی مرتبط می کنند. موفقیت ارتباط کمی ساختار-ویژگی را می توان توسط بینشی که به تخمین ساختاری خواص شیمیایی و امکان تخمین خواص ترکیبات شیمیایی جدید می انجامد توضیح داد. یک فرض اساسی در روش ارتباط کمی ساختار-ویژگی این است که همه خواص (فیزیکی- شیمیایی) از یک ماده شیمیایی به ساختار مولکولی مرتبط است. بر اساس این فرض، تحقیقات ارتباط کمی ساختار-ویژگی در دهه های اخیر توسعه بسیاری پیدا کرده اند.

1. Quantitative Structure-Property Relationship
2. Ordinary Least Squares
3. Principal Component Regression
4. Partial Least Squares
5. Artificial Neural Networks(ANN)
6. Support Vector Machine (SVM)

هدف از مطالعه حاضر، پیش بینی نقاط ذوب ترکیبات حلقوی نیتروآروماتیک با استفاده از توصیف کننده های شیمیایی و مکانیک کوانتومی و ماشین بردار پشتیبان است.

۲. بخش محاسباتی

۲-۱. سری داده ها

سری داده ها مربوط به نقاط ذوب ۶۰ ترکیب از نیتروآروماتیک های حلقوی است که توسط وانگ و همکارانش گزارش شده است [۱۸]. ساختار این ترکیبات در جدول ۱ آورده شده است. در ابتدا ترکیبات به صورت تصادفی به دو سری شامل سری آموزش و سری پیش بینی تقسیم شدند (جدول ۱). سری آموزش شامل ۳۵ ترکیب (۶۰٪ داده ها) و سری پیش بینی یا تست نیز شامل ۱۱ ترکیب (۲۰٪ داده ها) است. مقادیر نقاط ذوب به عنوان متغیر وابسته و توصیف کننده ها به عنوان متغیر مستقل انتخاب شدند.

۲-۳. محاسبه توصیف کننده ها

توصیف کننده ها مقادیر عددی هستند که خصوصیات مختلفی از مولکول را بیان می کنند. در حال حاضر تعداد زیادی توصیف کننده مولکولی وجود دارد که در مطالعات QSPR مورد استفاده قرار می گیرند. بعد از ارزیابی و پیدا کردن مناسب ترین آنها می توانند جهت پیش بینی خاصیت ترکیبات جدید بکار روند.

محاسبه توصیف کننده های الکترونی بوسیله بسته نرم افزاری Gaussian03W انجام شد. ساختار هندسی ۶۰ ترکیب حلقوی نیتروآروماتیک به روش DFT با تابع B3LYP و سری پایه 6-31G(d) بهینه شد. سپس چندین پارامتر ساختاری مرتبط (توصیف کننده) از نتایج محاسبات کوانتومی انتخاب شدند. که مشخصات این توصیف کننده ها در جدول ۲ آورده شده است. برای محاسبه بقیه توصیف کننده ها از نرم افزار دراگون استفاده گردید.

جدول ۲. مشخصات برخی از توصیف کننده های محاسبه شده.

توصیف کننده ها	نماد	اختصار	توصیف کننده ها	نماد	اختصار
توصیف کننده های مکانیک کوانتومی	Highest Occupied Molecular Orbital	HOMO	توصیف کننده های مکانیک کوانتومی	Hardness [$\eta=1/2 (HOMO+LUMO)$]	H
	Lowest Unoccupied Molecular Orbital	LUMO		Softness ($S=1/\eta$)	S
	difference between LUMO and HOMO	E GAP		Electro negativity [$\chi=-1/2 (HOMO-LUMO)$]	χ
	Molecular Polarizability	MP		Electrophilicity ($\omega=\chi^2/2\eta$)	Ω
ویژگی های شیمیایی	Molecule surface area (Approx)	SA (A)	ویژگی های شیمیایی	Partition Coefficient	Log P
	Molecule surface area (Grid)	SA (G)		Hydration Energy	HE
	Mass	M		Refractivity	REF
	Molecule volume	V			

۲-۴. ماشین بردار پشتیبان

ماشین بردار پشتیبان یکی از روش های یادگیری تحت نظارت است که هم برای دسته بندی و هم رگرسیون قابل استفاده است. این روش توسط وپنیک [۱۹] بر پایه تئوری یادگیری آماری بنا نهاده شده است. ماشین بردار پشتیبان روشی برای طبقه بندی دوتائی در فضای ویژگی های دلخواه است و از این روش مناسب برای مسائل پیش بینی به شمار می رود [۲۰]. ماشین بردار پشتیبان در اصل یک دسته بندی کننده دو کلاسه است که کلاس ها را توسط یک مرز خطی از هم جدا می کند. در این روش نزدیکترین نمونه ها به مرز تصمیم گیری را بردارهای پشتیبان می نامند. این بردارها معادله مرز تصمیم گیری را مشخص می کنند. الگوریتم های شبیه سازی معمولاً هوشمند کلاسیک مانند شبکه های عصبی مصنوعی، معمولاً قدر مطلق خطا یا مجموع مربعات خطای داده های آموزشی را حداقل می کنند، ولی مدل های SVM اصل حداقل سازی خطای ساختاری را به کار می گیرند [۲۱].

در یک مدل رگرسیونی SVM لازم است وابستگی تابعی متغیر وابسته y به مجموعه ای از متغیرهای مستقل X تخمین زده شود. فرض بر این است که مانند دیگر مسائل رگرسیونی، رابطه بین متغیرهای وابسته و مستقل توسط یک تابع معین f به علاوه یک مقدار اضافی نویز مشخص می شود.

$$y = f(x) + \text{Noise} \quad (1)$$

بنابراین موضوع اصلی، پیدا کردن فرم تابع f است که بتواند به صورت صحیح، موارد جدیدی را که SVM تاکنون تجربه نکرده است پیش بینی کند. این تابع به وسیله آموزش مدل SVM بر روی یک مجموعه داده به عنوان مجموعه آموزش که شامل فرآیندی به منظور بهینه سازی دائمی تابع خطا است. قابل دسترسی است. بر مبنای تعریف این تابع خطا، دو نمونه از مدل های SVM شناخته شده است که عبارتند از الف) مدل های رگرسیونی SVM نوع اول که مدل های SVM-v نیز نامیده می شوند و ب) مدل های رگرسیونی SVM نوع دوم که با نام SVM-ε شناخته شده هستند. در این مطالعه SVM-ε به دلیل کاربرد گسترده آن در مسائل رگرسیونی مورد استفاده قرار گرفت. برای این مدل، تابع خطا به صورت زیر تعریف می شود:

$$\frac{1}{2} W^T W + C \sum_{i=1}^N \delta_i + C \sum_{i=1}^N \delta_i^* \quad (2)$$

تابع خطای فوق لازم به ذکر است که با توجه به محدودیت های زیر کمینه گردد (۱۳):

$$W^T \phi(X_i) + b - y_i \leq \varepsilon + \delta_i^* \quad (3)$$

$$y_i - W^T \phi(x_i) - b \leq \varepsilon + \delta_i^* \\ \delta_i, \delta_i^* \geq 0$$

که در این رابطه C ثابت گنجایش، W بردار ضرایب، W^T ترانزپوز بردار ضرایب، δ_i, δ_i^* ضرایب کمبود، b ضریب ثابت، N الگوهای آموزش مدل و ϕ تابع کرنل است. اطلاعات کمی در مورد انتخاب تابع غیر خطی مناسب ϕ در دسترس می باشد. ماشین های بردار پشتیبان برای حل مسائل غیر خطی، ابعاد مسئله را از طریق توابع کرنل تغییر می دهند. انتخاب کرنل برای SVM که حجم داده های آموزشی و ابعاد بردار ویژگی بستگی دارد. به عبارت دیگر، بایستی با توجه به این پارامترها تابع کرنلی را انتخاب نمود که توانایی آموزش برای ورودی های مساله را داشته باشد. در عمل چهار نوع کرنل خطی، کرنل چهار جمله ای، کرنل تانژانت هیپربولینک و کرنل گوسی (RBF) به کار گرفته می شوند. در جدول ۳ معادلات برخی از کرنل های رایج ارائه شده اند.

تایع کرنل	نوع تایع
$K(x_i, x_j) = x_i^T \cdot x_j$	خطی
$K(x_i, x_j) = (\gamma x_i^T \cdot x_j + C)d$	چند جمله ای
$K(x_i, x_j) = \tanh(\gamma x_i^T \cdot x_j + C)$	تانزانته هیپربولیک
$K(x_i, x_j) = \exp(-\gamma x_i - x_j ^2)$	RBF

در نهایت، تایع تصمیم رگرسیون بردار پشتیبان غیر خطی به صورت معادله زیر خواهد بود که کنترل کننده میزان نوسان تایع گوسی و همچنین کنترل کننده نتایج پیش بینی و تعمیم دهنده مدل SVM است [۲۲].

$$f(x_i) = \sum_{i=1}^l (-\theta_i - \theta_i^*) K(x_i, x_j) + b \quad (۴)$$

۳. نتایج و بحث

۳-۱. مدل سازی با روش رگرسیون خطی چندگانه

برای ساختن مدلی که بیانگر ارتباط ساختاری ترکیبات مورد بررسی با نقطه ذوب (Mp) آنها باشد، از روش رگرسیون خطی چندگانه (MLR) استفاده شد. در ابتدا توصیف کننده های محاسبه شده، به عنوان متغیرهای مستقل و مقادیر Mp ترکیبات مورد نظر به عنوان متغیرهای وابسته، به عنوان ورودی به نرم افزار SPSS وارد شدند. در نهایت با استفاده از منوی آنالیز، گزینه ی رگرسیون خطی و روش مرحله ای انتخاب و نهایتاً چندین مدل مختلف بطور جداگانه به دست آمد، که با توجه به خصوصیات آماری آنها از جمله ضریب رگرسیون (R)، آماره F و خطای استاندارد و پس از رسم مقادیر R^2_{train} ، R^2_{test} ، $RMSE_{train}$ و $RMSE_{test}$ بر حسب تعداد توصیف کننده ها، بهترین مدل که دارای بیشترین مقدار R^2 و F و کمترین مقدار خطای استاندارد و شامل توصیف کننده های تا حد امکان قابل توجیه باشد، به عنوان مدل نهایی برای ارتباط Mp ترکیبات با ساختار آنها انتخاب شد. با این روش ۱۰ مدل بررسی شد که مدل هفتم با تعداد ۷ توصیف کننده به عنوان مناسب ترین آنها انتخاب و توسط روش MLR مدل سازی و مورد ارزیابی قرار گرفت. فهرست توصیف کننده های انتخاب شده توسط نرم افزار SPSS به همراه توصیف مختصری از آنها در جدول ۴ آورده شده است. جدول ۵ نیز مقادیر عددی توصیف کننده ها را نشان می دهد.

جدول ۴. توصیف کننده های انتخاب شده با SPSS و توصیف آنها.

توصیف کننده	نوع توصیف کننده	ضرایب	MF	VIF
MATS3e	2D autocorrelations	-298.37	0.076	1.099
Mor02m	3D-MoRSE descriptors	14.09	0.746	1.142

Mor32e	3D-MoRSE descriptors	96.74	-0.057	1.240
G2v	WHIM descriptors	92.80	0.090	1.079
H5m	GETAWAY descriptors	51.97	0.025	1.028
nHDon	Functional groups	28.10	0.085	1.032
E _{Gap}	Mechanic quantum descriptors	80.72	0.033	1.216
Constant	-	133.98	-	-

پس از انتخاب مناسب ترین توصیف کننده‌ها توسط روش مرحله ای با استفاده از SPSS، مرحله بعد، ایجاد مدل بین توصیف کننده های انتخاب شده و Mp ترکیبات است. از نرم افزار SPSS برای این منظور استفاده گردید و بین توصیف کننده‌ها و Mp ترکیبات سری آموزش با استفاده از روش MLR رابطه زیر به دست آمد:

$$Mp = 133.98 - 298.37 (\text{MATS3e}) + 14.09 (\text{Mor02m}) + 96.74 (\text{Mor32e}) + 92.80 (\text{G2v}) + 51.97 (\text{H5m}) + 28.10 (\text{nHDon}) + 80.72 (\text{E}_{\text{Gap}})$$

$$R^2_{\text{train}}=0.804, \quad F_{\text{train}}=23.54, \quad R^2_{\text{test}}=0.865, \quad F_{\text{test}}=4.022, \quad R^2_{\text{adj}}=0.770, \quad Q^2_{\text{LOO}}=0.721, \quad Q^2_{\text{LGO}}=0.594$$

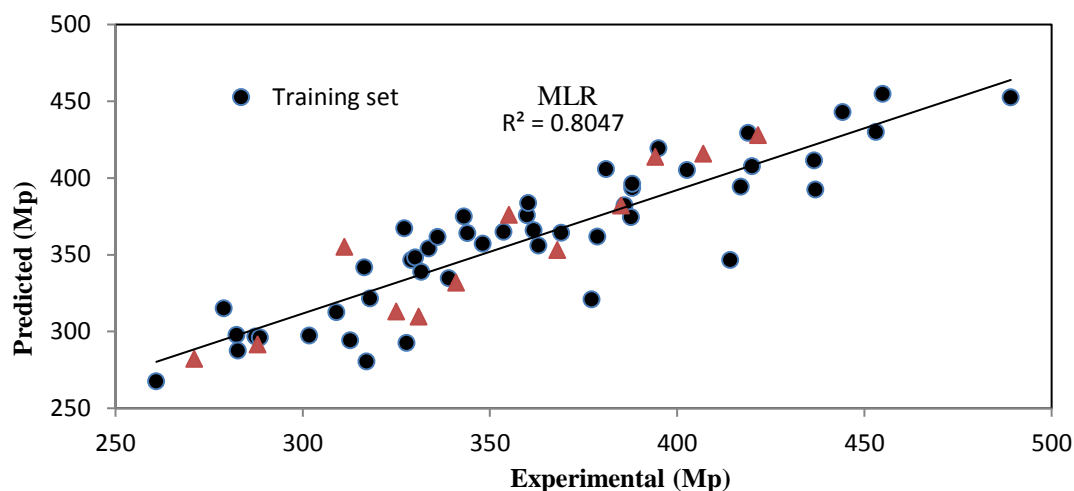
جدول ۵. مقادیر عددی توصیف کننده‌ها.

MATS3e	Mor02m	Mor32e	G2v	H5m	nHDon	E _{Gap}
0.054	9.902	-0.100	0.188	0.029	0	0.0097
-0.003	9.787	-0.144	0.416	0.000	0	0.2177
-0.009	10.866	-0.123	0.193	0.016	0	0.0150
-0.013	10.847	-0.233	0.188	0.026	0	0.0083
-0.007	10.596	-0.268	0.188	0.039	0	0.2193
-0.023	9.930	-0.169	0.197	0.000	0	0.1671
-0.032	10.648	-0.176	0.188	0.007	0	0.0366
0.165	11.283	0.051	0.204	0.000	1	0.2097
0.024	10.079	-0.238	0.193	0.024	1	0.0138
0.058	13.795	-0.171	0.171	0.040	1	0.0152
-0.002	10.906	-0.300	0.181	0.060	0	0.0149
-0.010	10.901	-0.178	0.204	0.000	1	0.0147
-0.157	11.760	-0.147	0.193	0.011	0	0.2027
0.032	11.382	-0.143	0.188	0.025	0	0.0356
-0.032	11.819	-0.163	0.191	0.489	0	0.1117
-0.016	11.671	-0.181	0.191	0.695	0	0.1077
0.005	10.740	-0.087	0.193	0.000	1	0.2165
-0.064	11.073	0.126	0.183	0.007	0	0.1919
-0.047	13.052	-0.219	0.217	0.029	1	0.0146
-0.062	12.241	-0.143	0.191	0.105	0	0.0043
-0.040	11.605	-0.120	0.191	0.018	2	0.0278
0.009	10.491	-0.064	0.200	0.013	2	0.1958

0.016	14.571	-0.256	0.175	0.046	1	0.0088
-0.168	12.427	-0.171	0.191	0.019	0	0.0429
-0.150	13.032	-0.080	0.188	0.057	0	0.0083
-0.093	12.475	-0.161	0.188	0.010	1	0.1948
0.026	11.770	-0.036	0.191	0.057	2	0.0053
-0.106	11.676	-0.101	0.200	0.000	0	0.2100
-0.109	11.886	-0.133	0.172	0.439	0	0.0552
0.005	11.456	-0.234	0.193	0.010	1	0.0387
0.107	16.522	-0.280	0.165	0.051	1	0.0982
-0.155	13.688	-0.211	0.197	0.094	1	0.0226
-0.091	10.046	-0.086	0.200	0.000	2	0.1627
-0.035	11.611	-0.089	0.242	0.834	0	0.1149
-0.155	12.509	-0.127	0.197	0.030	1	0.0162
-0.122	14.205	-0.016	0.236	0.000	0	0.0657
-0.203	13.742	-0.176	0.191	0.027	1	0.0114
-0.174	13.075	-0.120	0.177	0.487	0	0.0602
-0.084	11.265	-0.186	0.197	0.000	1	0.0042
-0.094	12.342	-0.176	0.197	0.388	1	0.1088
-0.067	10.847	-0.193	0.197	0.631	3	0.0715
-0.091	10.608	-0.040	0.310	0.005	2	0.1943
-0.107	12.290	-0.084	0.197	0.264	0	0.6009
-0.064	12.682	0.150	0.322	0.023	0	0.1861
-0.106	13.451	-0.059	1.000	0.000	0	0.0070
-0.113	12.275	-0.006	0.193	0.016	2	0.1848
-0.199	14.760	-0.249	0.188	0.063	2	0.0089
-0.294	13.793	-0.143	0.185	0.080	1	0.0113
0.003	10.351	-0.179	0.197	0.024	0	0.0107
0.011	10.442	-0.083	0.188	0.055	0	0.0163
-0.157	10.980	-0.142	0.193	0.035	0	0.1686
-0.023	10.491	-0.175	0.255	0.008	0	0.2129
-0.028	11.464	-0.166	0.193	0.024	0	0.0335
-0.078	12.261	-0.192	0.191	0.035	0	0.0087
-0.153	12.946	-0.145	0.185	0.113	0	0.0594
-0.158	10.123	-0.198	0.204	0.000	1	0.0166
-0.158	10.716	-0.140	0.204	0.001	1	0.2027
-0.226	12.559	0.005	0.193	0.000	0	0.2100
-0.202	11.999	-0.151	0.197	0.226	1	0.1107
-0.063	13.369	-0.012	0.191	0.017	2	0.1628

سپس از معادله به دست آمده برای پیش بینی Mp سری پیش بینی (تست) استفاده گردید. شکل ۱ نمودار مقادیر Mp محاسبه شده با کمک مدل SW-MLR برای مجموعه‌های آموزش و پیش‌بینی را برحسب مقادیر تجربی نشان می‌دهد.

با توجه به پارامترهای آماری، هر چند روش رگرسیون خطی چند گانه توانسته پیش بینی نسبتاً مناسبی را نشان دهد اما برای حصول نتایج بهتر، از ماشین بردار پشتیبان (SVM) برای پیش بینی Mp ترکیبات استفاده شد.



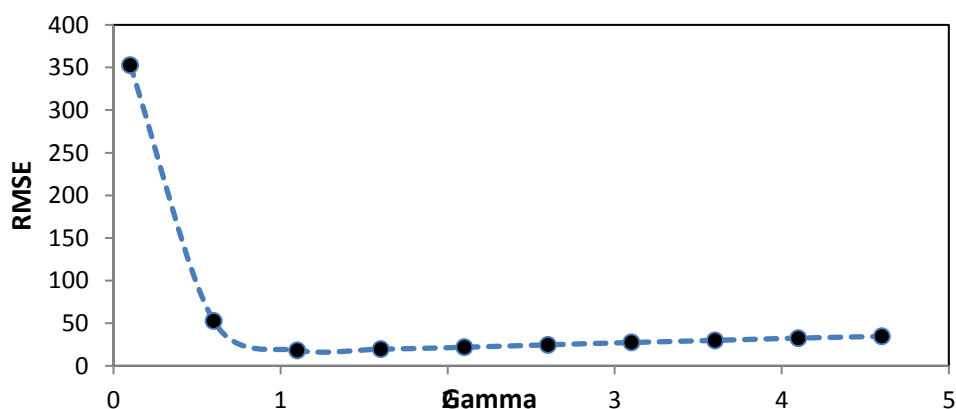
شکل ۱. نمودار مقادیر Mp محاسبه شده با کمک مدل SW-MLR برای مجموعه‌های آموزش و پیش‌بینی را بر حسب مقادیر تجربی.

۲-۳. مدل‌سازی با استفاده از ماشین بردار پشتیبان

در این مرحله توصیف‌کننده‌های انتخاب شده برای ایجاد مدل و پیش‌بینی فعالیت‌های بازدارندگی داروها توسط روش غیر خطی ماشین بردار پشتیبان مورد بررسی قرار گرفت. در این روش قبل از شروع مدل‌سازی لازم است یک سری از پارامترها بهینه شود. پارامترهایی که روی قدرت مدل‌سازی SVM تاثیر می‌گذارند عبارتند از:

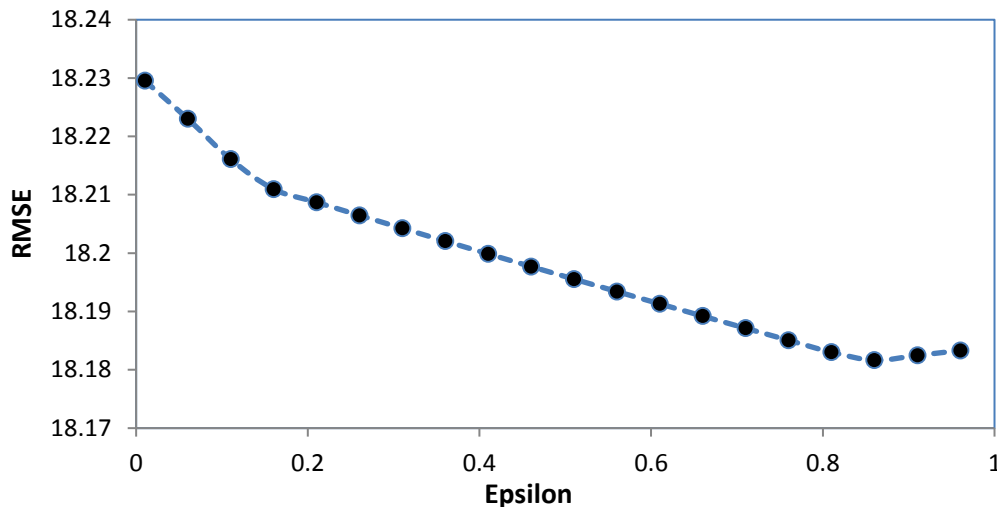
kernel function type, capacity parameter (C), ϵ -insensitive (ϵ), Gamma (eps)

در ابتدا باید نوع تابع مورد استفاده بهینه شود که نوع تابع، نحوه توزیع نمونه‌ها را در فضا مشخص می‌کند. RBF یکی از توابعی است که به طور معمول استفاده می‌شود و نتایج خوبی نیز می‌دهد در این پروژه از RBF به عنوان تابع برای SVM استفاده گردید. بعلاوه پارامتر متناظر با نوع تابع یعنی Gamma که روی تعداد بردارهای پشتیبان تاثیر می‌گذارد نیز باید بهینه شود. تعداد بردارهای پشتیبان بر زمان آموزش مدل تاثیر می‌گذارد طوری که افزایش مقدار Gamma و در نتیجه تعداد بردار پشتیبان می‌تواند به افزایش زمان آموزش و همچنین overfitting منجر شود. مقدار Gamma توانایی و قدرت SVM را در پیشگویی کنترل می‌کند. در شکل ۲ نمودار مقادیر متفاوت Gamma بر حسب RMSE نمایش داده شده است.



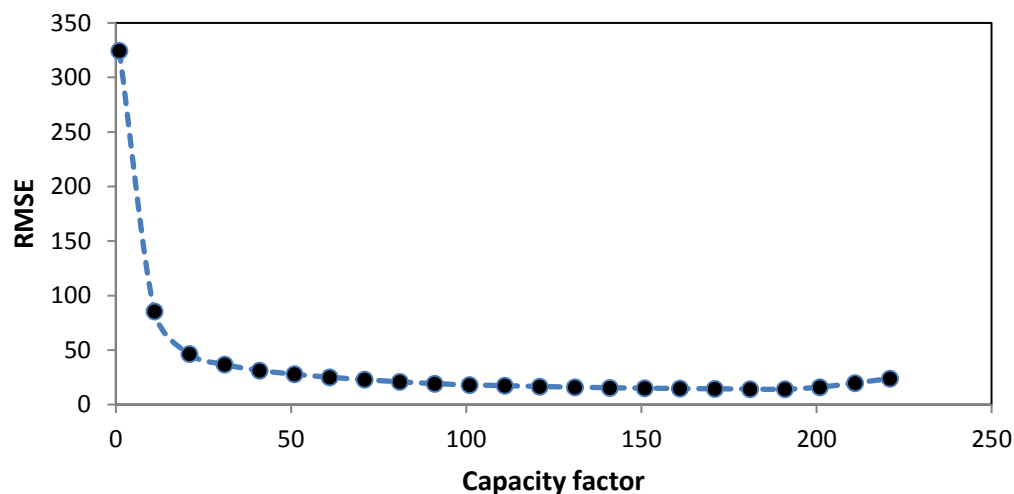
شکل ۲. نمودار تغییرات مقدار Gamma بر حسب مقدار RMSE برای سری آموزش

همانطوری که از این شکل ملاحظه می شود مقدار Gamma از ۰/۱ تا ۴/۶ تغییر می یابد و از نقطه ۱/۱ به بعد با افزایش مقدار Gamma RMSE افزایش می یابد. بنابراین مقدار ۱/۱ به عنوان نقطه بهینه برای Gamma انتخاب شد. پارامتر ϵ -insensitive parameter یا فاکتور حساسیت یکی دیگر از پارامترهایی است که باید بهینه شود. فاکتور حساسیت به نویزهای موجود در داده ها مربوط می شود که معمولاً ناشناس هستند. در شکل ۳ نمودار تغییرات ϵ (Epsilon) بر حسب RMSE نمایش داده شده است.



شکل ۳. نمودار تغییرات مقدار ϵ (Epsilon) بر حسب مقدار RMSE برای سری آموزش

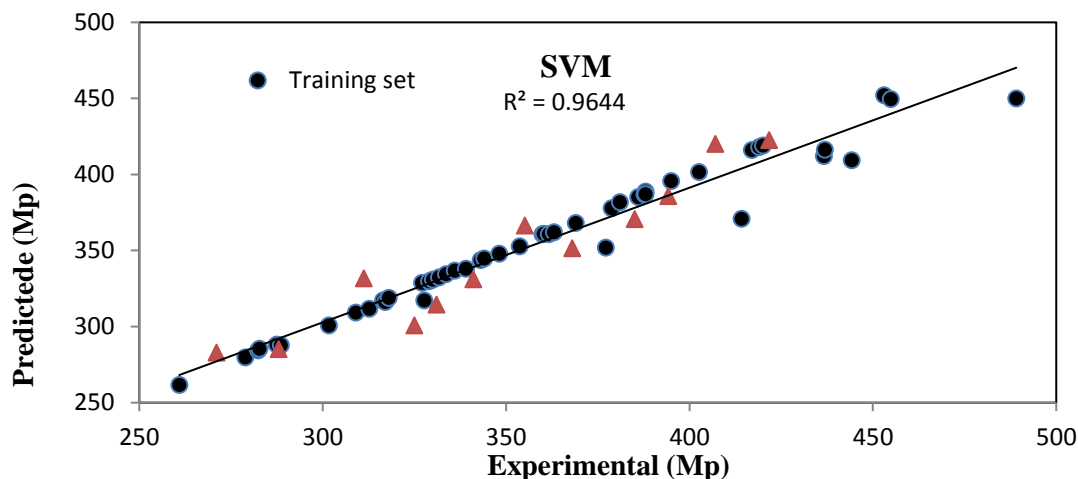
با تغییر مقدار Epsilon از ۰/۰۰۱ تا ۰/۶۹ مقدار بهینه ۰/۸۶ برای Epsilon انتخاب شد. و در نهایت پارامتر ظرفیت C باید بهینه شود اگر مقدار C پایین باشد یک پراکندگی در پیشگویی دیده خواهد شد و در بعضی اوقات با افزایش بیش از حد C، Overfitting رخ می دهد هر چند که مقدار زیادی C تاثیر چندانی روی پیشگویی ندارد ولی با این حال مقدار این پارامتر نیز باید بهینه شود. در شکل ۴ نمودار تغییرات Capacity parameter (C) بر حسب RMSE نمایش داده شده است.



شکل ۴. نمودار تغییرات مقدار Capacity parameter (C) بر حسب مقدار RMSE برای سری آموزش

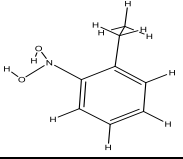
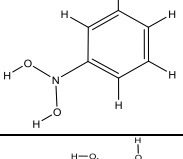
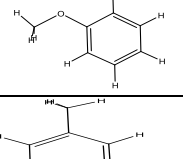
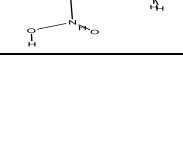
با تغییر مقدار Capacity parameter از ۱ تا ۲۲۱ مقدار بهینه ۱۹۱ برای Capacity parameter انتخاب شد. در مرحله آخر با استفاده از تمامی پارامترهای بهینه شده، مدل SVM ساخته شده و مقادیر فعالیتهای بازدارندگی ترکیبات دارویی پیش بینی شد. با استفاده از مدل SVM بهینه شده مقادیر فعالیتهای بازدارندگی ترکیبات مورد نظر در مجموعه آموزشی و پیش بینی مورد محاسبه قرار گرفته و در جدول ۱ نشان داده شده است

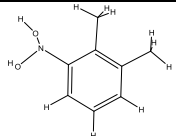
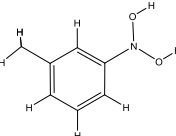
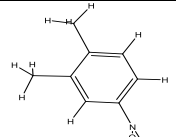
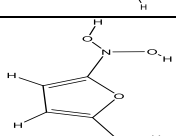
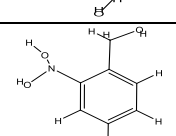
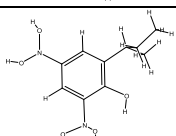
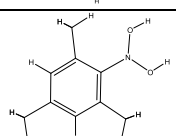
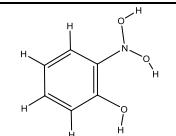
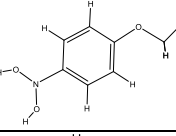
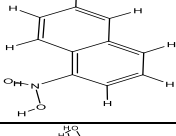
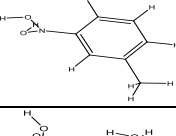
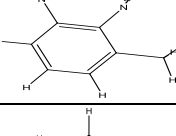
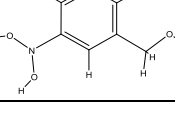
مقادیر فعالیتهای محاسبه شده و تجربی برای ترکیبات براساس مدل SVM در دو مجموعه آموزشی و تست در شکل ۵ آورده شده است در این شکل میزان نزدیکی داده ها به خط راست قدرت پیشگویی مدل را نشان می دهد.

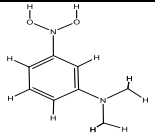
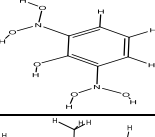
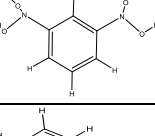
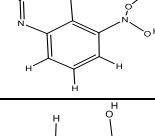
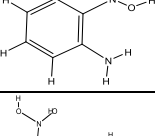
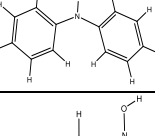
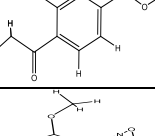
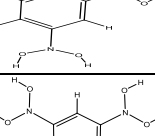
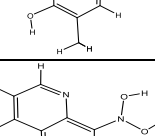
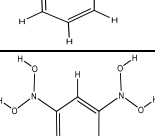
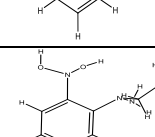
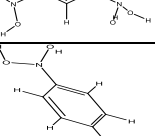
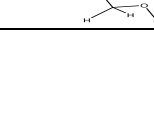


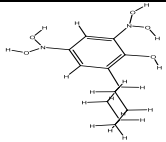
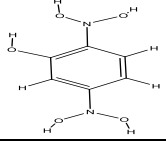
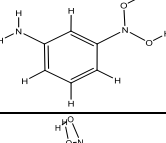
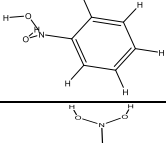
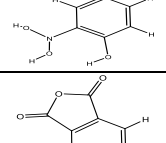
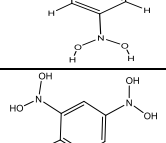
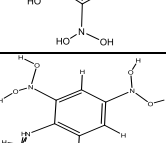
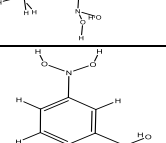
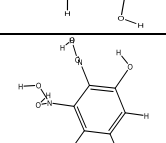
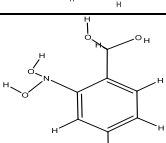
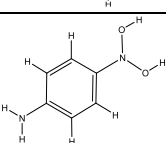
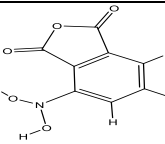

شکل ۵. مقادیر Mp محاسبه شده برای ترکیبات براساس مدل SVM در دو مجموعه آموزشی و تست بر حسب مقادیر تجربی

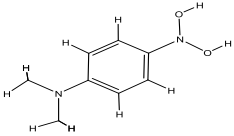
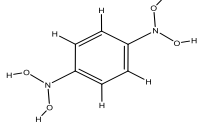
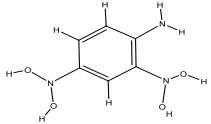
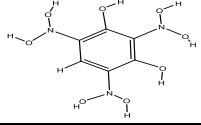
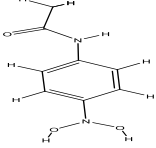
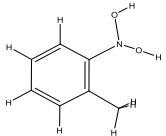
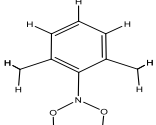
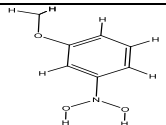
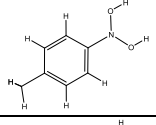
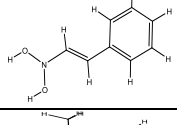
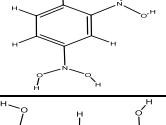
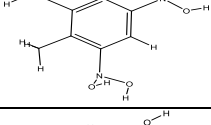
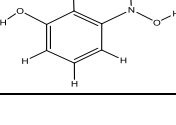
جدول ۱. مقادیر تجربی و محاسبه شده سمیت (pEC₅₀) برای ترکیبات مختلف برای مجموعه های آموزش، ارزیابی و پیش بینی در مدل SW-ANN همراه با مقادیر خطای نسبی

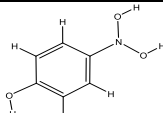
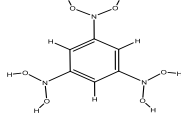
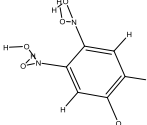
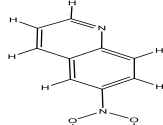
No.	Compound	Exp ^a (Mp)	MLR ^b (Mp)	SVM ^c (Mp)	Residual ^d (for SVM)
Training set					
1		260.90	267.48	267.48	6.58
2		278.90	315.05	315.05	36.15
3		282.35	297.85	297.85	15.50
4		282.68	287.65	287.65	4.97

5		287.40	296.65	296.65	9.25
6		288.59	296.21	296.21	7.62
7		301.70	297.33	297.33	-4.37
8		309.00	312.65	312.65	3.65
9		312.65	294.21	294.21	-18.44
10		316.42	341.82	341.82	25.40
11		317.00	280.37	280.37	-36.63
12		318.00	321.59	321.59	3.59
13		327.10	367.18	367.18	40.08
14		327.70	292.62	292.62	-35.08
15		329.00	346.48	346.48	17.48
16		330.00	348.26	348.26	18.26
17		331.65	338.92	338.92	7.27

18		333.65	354.15	354.15	20.50
19		336.00	361.68	361.68	25.68
20		339.00	334.68	334.68	-4.32
21		343.00	374.96	374.96	31.96
22		344.00	364.20	364.20	20.20
23		348.10	357.23	357.23	9.13
24		353.65	364.87	364.87	11.22
25		359.90	375.73	375.73	15.83
26		360.25	383.75	383.75	23.50
27		361.65	365.93	365.93	4.28
28		363.00	355.90	355.90	-7.10
29		369.00	364.38	364.38	-4.62
30		377.15	320.95	320.95	-56.20

31		378.70	361.79	361.79	-16.91
32		381.00	405.81	405.81	24.81
33		386.00	382.29	382.29	-3.71
34		387.70	374.52	374.52	-13.18
35		388.00	393.48	393.48	5.48
36		388.00	396.22	396.22	8.22
37		395.00	419.33	419.33	24.33
38		402.60	405.15	405.15	2.55
39		414.15	346.53	346.53	-67.62
40		417.00	394.26	394.26	-22.74
41		419.00	429.32	429.32	10.32
42		420.00	407.67	407.67	-12.33
43		436.60	411.49	411.49	-25.11

44		436.90	392.41	392.41	-44.49
45		444.20	442.83	442.83	-1.37
46		453.05	429.97	429.97	-23.08
47		454.90	454.91	454.91	0.01
48		489.10	452.59	452.59	-36.51
Test set					
1		271.00	282.04	282.04	11.04
2		288.00	291.45	291.45	3.45
3		311.15	355.17	355.17	44.02
4		325.00	313.03	313.03	-11.97
5		331.00	309.70	309.70	-21.30
6		341.00	331.72	331.72	-9.28
7		355.10	375.88	375.88	20.78
8		368.00	353.00	353.00	-15.00

9		385.00	382.05	382.05	-2.95
10		394.20	413.75	413.75	19.55
11		407.00	415.81	415.81	8.81
12		421.65	427.97	427.97	6.32

(a) مقادیر تجربی سمیت (pEC₅₀)

(b) مقادیر پیش بینی شده توسط رگرسیون خطی چندگانه

(c) مقادیر پیش بینی شده توسط ماشین بردار پشتیبان

(d) باقیمانده (اختلاف بین مقدار پیش بینی شده و مقدار واقعی برای مدل SVM)

۳-۳. ارزیابی اعتبار مدل های انتخاب شده

جهت اطمینان از اعتبار مدل بدست آمده، باید مدل را ارزیابی کرد. برای ارزیابی مدل از آماره های مختلفی نظیر F ، R^2 ، Q^2 ، $RMSE$ و آنالیز باقی مانده ها استفاده می شود. نتایج پارامترهای آماری بدست آمده در جدول ۶ آورده شده است. مقدار R^2 در واقع نشانگر میزان تطابق نتایج تجربی با مقادیر محاسباتی است یا به عبارتی دیگر معیاری از قدرت تشریح مدل رگرسیون می باشد.

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} \quad (5)$$

در معادله بالا y ، \hat{y} و \bar{y} به ترتیب نشان دهنده مقادیر تجربی، مقادیر محاسباتی و میانگین مقادیر تجربی هستند. برای مدل های MLR و SVM مقادیر R^2 به ترتیب 0.804 و 0.964 بدست آمد که نشان دهنده قدرت پیش بینی خوب مدل ها خصوصاً SVM می باشد. مقدار F تخمینی از اهمیت آماری معادله رگرسیون می باشد مزیت مقدار F بر R^2 این است که این پارامتر وابسته به درجه آزادی است و به طور واضح تر تحت تأثیر تعداد نمونه ها و تعداد توصیف کننده هایی است که در مدل دخیل هستند بنابراین ممکن است مدلی ساخته باشیم که R^2 خیلی بالا و مقدار F آن کم باشد این بدان معنی است که در مدل تعداد نمونه ها کم و در عوض تعداد توصیف کننده ها بالا بوده است

$$F = \frac{\sum_i (\hat{y}_i - \bar{y})^2 / p}{\sum_i (y_i - \hat{y}_i)^2 / (n - p - 1)} \quad (6)$$

در فرمول بالا n تعداد مولکولها و p نشانگر تعداد توصیف کننده های موجود در مدل به اضافه یک (عرض از مبدأ به عنوان یک درجه آزادی در نظر گرفته می شود) می باشد. در مدل بدست آمده برای MLR و SVM مقادیر F به ترتیب $23/54$ و $97/20$ بدست آمد. یکی دیگر از راه های ارزیابی مدل، پارامتر $RMSE$ می باشد که از رابطه زیر محاسبه خواهد شد:

$$RMSE = \frac{\sqrt{\sum_{i=1}^n (\hat{y}_i - y_i)^2}}{n} \quad (7)$$

پارامتر RMSE مقدار پراکندگی مقادیر تجربی حول خط رگرسیون را مشخص می‌کند. پایین بودن مقدار RMSE نشان دهنده اعتبار بالای مدل می‌باشد. مقادیر RMSE برای سری های آموزش و پیش بینی در هر دو مدل در جدول ۶ آورده شده است.

جدول ۶. پارامترهای آماری بدست آمده برای مدل های SVM و MLR

	Training set			Test set		
	R ²	RMSE	F	R ²	RMSE	F
MLR	0.804	23.299	23.549	0.865	18.073	4.022
SVM	0.964	11.566	97.207	0.910	14.040	6.093

۴. نتیجه گیری

در این مطالعه، رگرسیون QSPR برای پیش بینی نقاط ذوب ترکیبات پرانرژی نیتروآروماتیک های حلقوی بررسی شده است. در ابتدا توصیف کننده های الکترونی و مکانیک کوانتومی و شیمیایی برای هر ترکیب بدست آمد. سپس رابطه ای بین مناسب ترین توصیف کننده ها و نقاط ذوب ترکیبات با استفاده از روش های SVM و MLR بدست آمد. نتایج نشان داد که هر دو مدل SVM و MLR نتایج قابل قبولی ارائه داده است اگر چه مقایسه این دو روش نشان از برتری روش SVM نسبت به روش MLR دارد. بنابراین می‌توان نتیجه گرفت که از توصیف کننده های الکترونی و مکانیک کوانتومی نیز می‌توان جهت پیش بینی نقاط ذوب ترکیبات فوق استفاده کرد. از طرف دیگر از این روش ها می‌توان جهت پیش بینی سایر خواص فیزیکوشیمیایی ترکیبات نیز استفاده نمود.

۵. مراجع

- [1] Hamadani, M., Keshavarz, M.H., Nazari, B. and Mohebbi, M., Reliable method for safety assessment of melting points of energetic compounds. *Process Safety and Environmental Protection*, 103 (2016) 10-22.
- [2] Pagoria, P.F., Lee, G.S., Mitchell, A.R. and Schmidt, R.D., A review of energetic materials synthesis. *Thermochimica Acta*, 384(1) (2002) 187-204.
- [3] Rice, B.M., Pai, S.V. and Hare, J., Predicting heats of formation of energetic materials using quantum mechanical calculations. *Combustion and flame*, 118(3) (1999) 445-458.
- [4] Shahlaie, M., Fassihi, A., Pourhossein, A. and Arkan, E., Statistically validated QSAR study of some antagonists of the human CCR5 receptor using least square support vector machine based on the genetic algorithm and factor analysis. *Medicinal Chemistry Research*, 22(3) (2013) 1399-1414.
- [5] Ma, S., Lv, M., Deng, F., Zhang, X., Zhai, H. and Lv, W., Predicting the ecotoxicity of ionic liquids towards *Vibrio fischeri* using genetic function approximation and least squares support vector machine. *Journal of hazardous materials*, 283 (2015) 591-598.
- [6] Prasoona, R.K., Jyoti, A., Mukesh, Y., Nishant, S., Anuraj, N.S. and Shobha, J., Optimization of Gaussian Kernel Function in Support Vector Machine aided QSAR studies of C-aryl glucoside SGLT2 inhibitors. *Interdisciplinary Sciences: Computational Life Sciences*, 5(1) (2013) 45-52.

- [7] Nekoei, M., Mohammadhosseini, M. and Pournasheer, E., QSAR study of VEGFR-2 inhibitors by using genetic algorithm-multiple linear regressions (GA-MLR) and genetic algorithm-support vector machine (GA-SVM): a comparative approach. *Medicinal Chemistry Research*, 24(7) (2015) 3037-3046.
- [8] Pournasheer, E., Aalizadeh, R., Ganjali, M.R. and Norouzi, P., Prediction of Superoxide Quenching Activity of Fullerene (C60) Derivatives by Genetic Algorithm-Support Vector Machine. *Fullerenes, Nanotubes and Carbon Nanostructures*, 23(4) (2015) 290-299.
- [9] Zhao, H., Zhang, X., Ji, L., Hu, H. and Li, Q., Quantitative structure-activity relationship model for amino acids as corrosion inhibitors based on the support vector machine and molecular design. *Corrosion Science*, 83 (2014) 261-271.
- [10] Pournasheer, E., Beheshti, A., Vahdani, S., Nekoei, M., Danandeh, M., Abbasghorbani, M. and Ganjali, M.R., Simple QSPR Modeling for Prediction of the GC Retention Indices of Essential Oil Compounds. *Journal of Essential Oil Bearing Plants*, 18(6) (2015) 1298-1309.
- [11] Jalali, A., Nekoei, M. and Mohammadhosseini, M., Novel QSPR Study on the Melting Points of a Broad Set of Drug-Like Compounds Using the Genetic Algorithm Feature Selection Approach Combined With Multiple Linear Regression and Support Vector Machine. *Journal of Chemical Health Risks*, 6(1) (2016).
- [12] Nekoei, M., A Robust Quantitative Structure-Property Relationship-Based Model for Estimation of Refractivity Indices of 101 Common Paraffin Derivatives Based on Their Molecular Structures. *Journal of Chemical Health Risks*, (2016).
- [13] Mohammadhosseini, M., Deeb, O., Alavi-Gharabagh, A. and Nekoei, M., Exploring novel QSRRs for simulation of gas chromatographic retention indices of diverse sets of terpenoids in Pistacia lentiscus L. essential oil using stepwise and genetic algorithm multiple linear regressions. *Analytical Chemistry Letters*, 2(2) (2012) 80-102.
- [14] Adimi, M., Salimi, M., Nekoei, M., Pournasheer, E. and Beheshti, A., A quantitative structure-activity relationship study on histamine receptor antagonists using the genetic algorithm-multi-parameter linear regression method. *J Serb Chem Soc*, 77(5) (2012) 639-650.
- [15] Rahimi, M. and Nekoei, M., Quantitative Structure-Property Relationship Study for Prediction of Flash Point of Some Organic Compounds Based On SW-MLR Method. *Analytical Chemistry Letters*, 3(4) (2013) 278-286.
- [16] Nekoei, M., Goudarzi, N., Nekoei, S. and Mohammadhosseini, M., QSAR Study of Arylsulfonylpiperazine Inhibitors of 11 β -HSD1 by GA-MLR, GA-PLS and GA-ANN. *Analytical Chemistry Letters*, 4(1) (2014) 14-28.
- [17] Beheshti, A., Pournasheer, E., Nekoei, M. and Banaei, A., Quantitative Structure-Activity Relationship Study of Amino Acid Derivatives as Histone Deacetylase Inhibitors using the Genetic Algorithm-Multiple Linear Regression. *Analytical Chemistry Letters*, 2(1) (2012) 33-43.
- [18] Wang, D., Yuan, Y., Duan, S., Liu, R., Gu, S., Zhao, S., Liu, L. and Xu, J., QSPR study on melting point of carbocyclic nitroaromatic compounds by multiple linear regression and artificial neural network. *Chemometrics and Intelligent Laboratory Systems*, 143 (2015) 7-15.
- [19] Ma, S., Lv, M., Deng, F., Zhang, X., Zhai, H. and Lv, W., Predicting the ecotoxicity of ionic liquids towards *Vibrio fischeri* using genetic function approximation and least squares support vector machine. *Journal of hazardous materials*, 283 (2015) 591-598.
- [20] Golmohammadi, H., Dashtbozorgi, Z. and Vander Heyden, Y., Support vector regression based QSPR for the prediction of retention time of peptides in reversed-phase liquid chromatography. *Chromatographia*, 78(1-2) (2015) 7-19.
- [21] Ghanbari, A., Application of support vector machine in QSAR study of triazolyl thiophenes as cyclin dependent kinase-5 inhibitors for their anti-alzheimer activity. *Indian Journal of Chemical Technology (IJCT)*, 23(1) (2016) 9-21.
- [22] Golmohammadi, H. and Dashtbozorgi, Z., QSPR studies for predicting polarity parameter of organic compounds in methanol using support vector machine and enhanced replacement method. *SAR and QSAR in Environmental Research*, 27(12) (2016) 977-997.

