

فهرست

فصل اول.....	۱
مقدمه‌ای بر کاربرد کمومتری‌کس در مطالعات روابط کمی بین ساختار و فعالیت / خاصیت.....	۱
۱.۱ مقدمه.....	۱
۱.۱.۱ شیمی.....	۱
۱.۱.۲ کمومتری‌کس.....	۲
۱.۱.۳ طراحی دارو.....	۵
۱.۲ کمومتری‌کس در QSAR/QSPR.....	۵
۱.۳ داده‌های مورد استفاده در این مطالعه.....	۸
۱.۴ روش‌های مدلسازی.....	۹
۱.۴.۱ رگرسیون خطی چندتابی:.....	۱۰
۱.۴.۲ رگرسیون مرزی:.....	۱۲
۱.۴.۳ روش‌های رگرسیون چندمتغیره مبتنی بر متغیرهای پنهان.....	۱۲
۱.۴.۴ تعیین تعداد بهینه‌ی متغیرهای پنهان.....	۱۶
۱.۵ وابستگی در بین متغیرهای مستقل و نتایج آن بر روی مدل.....	۱۷
۱.۶ عدم قطعیت در ضرایب رگرسیون.....	۱۹
۱.۶.۱ روش‌های اندازه‌گیری عدم قطعیت.....	۲۰
۱.۶.۱.۱ انتشار خطای تصادفی.....	۲۰
۱.۶.۱.۲ نمونه‌برداری مجدد.....	۲۲
۱.۶.۱.۳ خودراه‌انداز.....	۲۲
۱.۶.۱.۴ خودراه‌انداز نمونه‌ها.....	۲۲
۱.۶.۱.۵ خودراه‌انداز باقیمانده‌ها.....	۲۴
۱.۶.۱.۶ جک‌نایف.....	۲۵
۱.۷ روش‌های پیش‌پردازش ماتریس داده.....	۲۶
۱.۸ ضرورت انتخاب متغیر.....	۲۷
۱.۸.۱ مزایای انتخاب متغیر.....	۳۰
۱.۸.۲ طبقه‌بندی کلی روش‌های تخمین و یا انتخاب متغیر.....	۳۱
۱.۸.۳ انتخاب متغیر توسط رویه‌های تکاملی.....	۳۲

۳۳.....	۱.۸.۲ انتخاب متغیر توسط بردارهای حاوی اطلاعات مفید.....
۳۴.....	۱.۸.۲.۱ وابستگی با متغیر معیار.....
۳۴.....	۱.۸.۲.۲ ضرایب رگرسیون.....
۳۴.....	۱.۸.۲.۳ علامت خالص آنالیت.....
۳۵.....	۱.۸.۲.۴ بردار ظرفیت.....
۳۶.....	۱.۸.۲.۵ قدرت نفوذ.....
۳۷.....	۱.۸.۲.۶ نسبت گزینش پذیری.....
۳۹.....	۱.۸.۲.۷ اهمیت متغیر در تصویرسازی.....
۴۰.....	۱.۸.۲.۳ چگونگی کاربرد عدم قطعیت در انتخاب متغیر.....
۴۲.....	۱.۸.۲.۴ روش‌های مبتنی حذف و یا کاهش وابستگی در بین متغیرهای مستقل.....
۴۳.....	۱.۸.۲.۴.۱ عمودسازی گرم-اشمیت.....
۴۵.....	۱.۸.۲.۴.۲ رویه‌ی تصویرسازی متناوب.....
۴۷.....	۱.۹ اعتبارسنجی مدل.....
۴۷.....	۱.۹.۱ روش‌های اعتبارسنجی.....
۴۷.....	۱.۹.۱.۱ اعتبارسنجی تقاطعی و خودراه‌انداز.....
۴۷.....	۱.۹.۱.۲ آزمایش تصادفی کردن γ
۴۸.....	۱.۹.۱.۳ اعتبارسنجی خارجی.....
۴۹.....	۱.۹.۲ روابط پرکاربرد در سنجش اعتبار مدل.....
۵۰.....	فصل دوم
۵۰.....	مروری بر پژوهش‌های گذشته
۵۰.....	۱.۲ نگاهی به میزان اهمیت روش‌های انتخاب متغیر.....
۵۱.....	۲.۲ نگاهی بر پژوهش‌های صورت گرفته در مورد داده‌های همبسته.....
۵۳.....	۳.۲ نگاهی بر پژوهش‌های صورت گرفته به منظور کاهش و حذف همبستگی.....
۵۴.....	۴.۲ نگاهی بر پژوهش‌های صورت گرفته در زمینه‌ی انتخاب متغیر به وسیله تخمین عدم قطعیت.....
۵۶.....	فصل سوم
۵۶.....	نتایج و بحث
۵۶.....	بخش اول
۵۶.....	۱.۳ مقدمه.....
۵۷.....	۱.۳.۱ بهبود عملکرد روش SPA به وسیله‌ی تاثیر هر متغیر در پیش‌بینی مدل.....

۵۸.....	۱.۱.۳ روش GSO.....
۶۰.....	۱.۱.۳ روش SPA.....
۶۳.....	۱.۱.۳ عملکرد روش CWSPA در انتخاب متغیر.....
۶۵.....	۱.۱.۳ ارزیابی عملکرد روش QWSPA.....
۷۳.....	۱.۱.۳ معبرسازی تقاطعی مدل با استفاده از نمونه‌های خارجی آزمایش.....
۷۵.....	بخش دوم.....
۷۵.....	۲.۳ مقدمه.....
۷۵.....	۲.۳.۱ انتخاب غیرنظارتی متغیرهای حاوی اطلاعات مفید با استفاده از رویه‌ی جک‌نایف در داده‌ی HIV.....
۷۸.....	۲.۳.۱.۱ اعمال رویه‌ی جک‌نایف - PLS بر روی مجموعه‌ی عمودشده و بردارهای امتیاز.....
۸۳.....	۲.۳.۱.۲ تعیین تعداد بهینه‌ی بردارهای امتیاز.....
۸۶.....	۲.۳.۱.۳ تعیین مقدار SS.....
۸۷.....	۲.۳.۱.۴ حذف متغیرهای بی‌ارزش.....
۸۹.....	۲.۳.۱.۵ معبرسازی تقاطعی مدل با استفاده از نمونه‌های خارجی آزمایش.....
۹۲.....	مراجع.....

فهرست شکل‌ها

- شکل (۱-۱): پروفیسور سوانته ولد ۳
- شکل (۲-۱): شمای عمومی خط مشی *QSAR/QSPR* ۷
- شکل (۳-۱): شمای کلی رابطه‌ی بین ماتریس‌های X و Y در رویه‌ی *PLS* ۱۴
- شکل (۴-۱): میزان بیان متغیر وابسته (Y) در رگرسیون توسط دو متغیر مستقل X_1 و X_2 ، در حالت‌هایی که وابستگی بین دو متغیر مستقل (قسمت هاشورخورده) الف) هیچ، ب) کم و ج) زیاد است. ۱۸
- شکل (۵-۱): شمای نمونه‌برداری مجدد به روش خودراه‌انداز ۲۳
- شکل (۶-۱): شمای نمونه‌برداری مجدد به روش خودراه‌انداز باقیمانده‌ها ۲۵
- شکل (۷-۱): شمایی از نمونه‌برداری مجدد به روش جک‌نایف ۲۶
- شکل (۸-۱): چگونگی انتخاب متغیر به وسیله‌ی نمونه‌برداری مجدد (جک‌نایف) و آزمایش t ۴۱
- شکل (۹-۱): روند انتخاب متغیر در *GSO* و *SPA* برای داده‌ای با ۳ سطر و ۵ ستون نشان داده شده‌است. انتخاب متغیرها با شروع از X_2 شروع میشود و تصویرسازی در مراحل ۱ و ۲ نشان می‌دهد که به ترتیب متغیرهای X_3 و X_4 انتخاب شده‌اند. متغیرهای انتخاب شده توسط روش *GSO* و *SPA* برای این داده به ترتیب (X_2 و Px_3 و Ppx_5) و (X_2 و X_3 و X_5) هستند. ۴۵
- شکل ۱-۳: نمایش میزان همبستگی در بین متغیرهای داده‌ی *HIV*، الف) داده‌ی اصلی ب) ماتریس بدست‌آمده پس از اعمال *GSO* ۵۸
- شکل (۲-۳): الف) عدم یکسانی بین مجموعه‌های عمودشده‌ی درجه‌بندی و آزمایش و ب) راهکار پیشنهادی به منظور یکنواخت-سازی فضای دو مجموعه‌ی عمودشده‌ی درجه‌بندی و پیش‌بینی ۵۹
- شکل (۳-۳): نمایش میزان همبستگی در بین متغیرهای داده‌ی *HIV*، الف) داده‌ی اصلی ب) ماتریس بدست‌آمده پس از اعمال *SPA* ۶۰
- شکل (۴-۳): تغییرات قدرت پیش‌بینی مدل برای مجموعه‌های درجه‌بندی و اعتبارسنجی برای داده‌های الف) *HIV* ب) *GABA* ج) *FLUOR* ۶۲
- شکل (۵-۳): عملکرد روش *CWSPA* در انتخاب متغیر با توجه به توان بردار همبستگی (از بالا به پایین به ترتیب ۱ و ۳ و ۵) ۶۴
- شکل (۶-۳): نمودار بافت‌نمای مربوط به همبستگی موجود در بین متغیرهای انتخابی به روش *CWSPA* که توان بردار همبستگی در الف) ۱، ب) ۳ و ج) ۵ ۶۶
- شکل (۷-۳): تغییرات قدرت پیش‌بینی مدل با افزایش متغیرهای انتخاب شده در روش *QWSPA* و در حضور نمونه‌برداری مجدد برای مجموعه‌های درجه‌بندی و اعتبارسنجی برای داده‌های الف) *HIV* ب) *GABA* ج) *FLUOR* ۶۸

شکل (۳-۸): تغییرات قدرت پیش‌بینی مدل با افزایش متغیرهای انتخاب شده در روش *QWSPA* و در حضور مجموعه‌های درجه‌بندی و اعتبارسنجی ثابت برای داده‌های الف) *HIV* (ب) *GABA* (ج) *FLUOR* ۶۹

شکل (۳-۹): تغییرات قدرت پیش‌بینی مدل با افزایش متغیرهای انتخاب شده با رابطه‌ی (۳-۲) در روش *QWSPA* و در حضور مجموعه‌های درجه‌بندی و اعتبارسنجی ثابت برای داده‌های الف) *HIV* (ب) *GABA* (ج) *FLUOR* ۷۱

شکل (۳-۱۰): نمودار بافت‌نمای مربوط به همبستگی موجود در بین متغیرهای انتخابی به روش *QWSPA* برای سه داده‌ی الف) *HIV* (ب) *GABA* (ج) *FLUOR* ۷۲

شکل (۳-۱۱): نمودار جعبه‌ای برای ۱۰۰ بار تکرار مدل‌سازی و پیش‌بینی مجموعه‌ی خارجی آزمایش به وسیله‌ی ۶ متغیر ابتدایی انتخاب‌شده به وسیله‌ی روش *QWSPA* برای داده‌های الف) *HIV* (ب) *GABA* (ج) *FLUOR* ۷۳

شکل (۳-۱۲): *p-value* همه‌ی متغیرهای عمودشده‌ای که از فرایند جک‌نایف در مدل‌سازی بردار اول و سوم امتیاز بدست‌آمده است. *SS* برابر ۵۸ است. ۷۹

شکل (۳-۱۳): حدود اطمینان برای ضرایب رگرسیون متغیرهای عمودشده که توسط مدل‌کردن بردارهای الف) اول، ب) دوم، ج) سوم و د) چهارم در روش جک‌نایف و به صورت جداگانه بدست آمده‌است. ۸۰

شکل (۳-۱۴): الف) $Q^2_{TOT}(f3)$ (ب) Q^2_{EXT} و ج) تعداد متغیرهای معتبر برای هر بردار امتیاز که از اعمال روش جک‌نایف بر روی مجموعه‌ی متغیرهای عمودشده و هر یک از بردارهای ماتریس امتیاز (به صورت جداگانه) بدست آمده‌است. ۸۲

شکل (۳-۱۵): پیامدهای تعیین تعداد بهینه‌ی بردارهای امتیاز برای داده‌ی *HIV* که با روش‌های الف) تک حذفی ب) بردار ویژه ۸۴

شکل (۳-۱۶): تعداد متغیرهای کل معتبر بر حسب تعداد بردارهای امتیاز بکاررفته در فرایند جک‌نایف در مقادیر مختلف *SS* از ۳۲ (پایین‌ترین خط) تا ۷۰ (بالا‌ترین خط). ۸۵

شکل (۳-۱۷): توانایی پیش‌بینی بردارهای امتیاز توسط متغیرهای عمودشده‌ی معتبر با افزایش *SS* ۸۶

شکل (۳-۱۸): تفاوت در پیش‌بینی برای حضور و عدم‌حضور متغیر مورد نظر در ۱۰ مجموعه‌ی تصادفی از متغیرهای انتخابی برای بردارهای امتیاز الف) اول، ب) دوم، ج) سوم و د) چهارم ۸۸

شکل (۳-۱۹): توانایی پیش‌بینی فعالیت‌های ضد *HIV* با استفاده از مجموعه‌ی بدست‌آمده از روش غیرنظارتی جک‌نایف. نمودار جعبه‌ای از ۱۰۰ بار تکرار مدل‌سازی و پیش‌بینی مجموعه‌ی خارجی بدست آمده‌است. ۹۰

فهرست جدول‌ها

- جدول (۱-۱) : تاثیر وابستگی بر روی ضرایب رگرسیون در روشهای PLS و MLR ۱۹
- جدول (۱-۳) : اندیس بردارهای آغازگر برای دادههای HIV و $GABA$ و $FLUOR$ ۶۱
- جدول (۲-۳) : میزان پیشبینی مجموعههای خارجی برای دادههای HIV و $GABA$ و $FLUOR$ با استفاده از ۶ متغیر ابتدایی انتخابشده از روش $QWSPA$ ۷۴
- جدول (۳-۳) : متغیرهای معتبر و شاخصهای کیفیت مدل ($Q^2_{TOT(F3)}$ و R^2) که از بکاربردن روش جکنایف بر روی مجموعهی متغیرهای عمودشده و ماتریس امتیازها بدست آمدهاست. ۸۱
- جدول (۴-۳) : متغیرهای معتبر تجمعی بدستآمده از اعمال رویهی جکنایف بر روی چهار بردار امتیاز هنگامی که SS برابر ۶۰ است. ۸۷
- جدول (۵-۳) : متغیرهای انتخابی از روش جکنایف برای چهار بردار امتیاز ابتدایی و با توجه به تاثیر آنها بر روی پیشبینی ۸۹

فصل اول

مقدمه‌ای بر کاربرد کمومتری‌کس در مطالعات روابط کمی بین ساختار و فعالیت / خاصیت

۱.۱ مقدمه

۱.۱.۱ شیمی

شیمی^۱ شاخه‌ای از علوم فیزیکی^۲ است که وظیفه‌ی مطالعه‌ی ترکیب، خواص و رفتار ماده را به عهده دارد. گاهاً شیمی را به این علت که رابطه‌ای برای فیزیک و دیگر علوم طبیعی مانند زمین‌شناسی و زیست-شناسی است به نام "علم مرکزی"^۳ نیز می‌خوانند. اگرچه شیمی شاخه‌ای از علوم فیزیکی است، اما از علم فیزیک جداست [۱].

اگر شخصی در پی ریشه‌یابی کلمه‌ی شیمی باشد پس از سفرهای فراوان در زمان‌های قدیم و چالش کردن با تمدن‌هایی همچون اعراب، یونانیان، مصریان، چینی‌ها و یهودیان باستان و همچنین برخورد با کلماتی از قبیل آلکمی^۴، آلکیمی^۵ و کیمی^۶ باز هم به درستی به منشاء این کلمه نخواهد رسید. با تمام این قضایا در قرن ۱۸ میلادی واژه‌ی شیمی برای علم امروزی انتخاب شد.

¹ Chemistry

² Physical Sciences

³ Central Science

⁴ alchemy

⁵ Al-kimiya

⁶ kimi

در قرن ۲۰ میلادی و با گسترش علوم و نفوذ و همپوشانی آنها در برخی از موارد باعث شد که شاخه‌های جدیدی از علم و به نام‌های ترکیبی به علوم افزوده شود. در شیمی این شاخه‌ها با ترکیب شدن سایر واژه‌ها با دگر شکل‌هایی از واژه‌ی شیمی مانند کمو^۱، کمی^۲ و کم^۳، به شکل کموسفر^۴، کموتراپی^۵، کمیسورپشن^۶ و کمینفورماتیکس^۷، بوجود آمدند.

در این زمان بود که چندین رشته‌ی علمی با نام‌های ترکیبی که انتهای آنها به متری (مانند ژئومتری^۸) و متریکس (مانند بایومتریکس^۹) که از یک واژه‌ی یونانی به معنای اندازه‌گیری گرفته شده بود. از واژه‌های مشابه که برای نام‌گذاری علوم جدید استفاده شده، می‌توان به سایکومتری^{۱۰} و سایکومتریکس^{۱۱} در رشته‌ی روانشناسی و اکونومتری^{۱۲} و اکونومتریکس^{۱۳} در اقتصاد اشاره کرد.

۱. ۱. ۲. کمومتریکس

در سال ۱۹۷۱ میلادی بود که برای اولین بار واژه‌ی کمومتری^{۱۴} در زبان سوئدی و کمومتریکس^{۱۵} در زبان انگلیسی به عنوان ترکیبی از شیمی و اندازه‌گیری توسط یک شیمیدان آلی به نام پروفیسور سوانته ولد^{۱۶} ارائه شد. پروفیسور سوانته ولد کمومتریکس را در سال ۱۹۷۴ به عنوان هنری برای بیرون کشیدن اطلاعات موردنظر شیمیایی از داده‌هایی که در آزمایش‌های شیمیایی تهیه شده‌اند، تعریف کرد. ایشان در بیست سال پس از این تعریف، در مقاله‌ای [۲] تعریف کمومتریکس را به این شکل اصلاح کردند:

¹ Chemo

² Chemi

³ Chem

⁴ Chemosphere

⁵ Chemotherapy

⁶ Chemisorption

⁷ Cheminformatics

⁸ Geometry

⁹ Biometrics

¹⁰ Psychometry

¹¹ Psychometrics

¹² Econometry

¹³ Econometrics

¹⁴ Kemometri

¹⁵ Chemometrics

¹⁶ Svante Wold

دادن آن‌ها با علم شیمی دستورکار بهینه‌ای را برای بدست آوردن داده‌های مورد نیاز و به مفیدترین شکل آن ارائه دهد. معمولاً این مورد با استفاده از روش‌های طراحی آزمایش^۱ انجام می‌شود [۴].

از آنجاییکه معمولاً نمی‌توان از داده‌های به اصطلاح خام بدست آمده از آزمایش‌ها اطلاعات موردنظر را در مورد سیستم مورد بررسی بدست آورد، وظیفه‌ی بعدی و اصلی کمومتریکس در اینجا بیان می‌شود. دانشمندان عرصه‌ی کمومتریکس با استناد به قوانین، اصول و روش‌های آمار و ریاضی، در پی بهبود روش‌های موجود و طرح روش‌های جدیدتر به تحلیل داده‌های خام می‌پردازند. سپس شیمیدان می‌تواند با استفاده‌ی از این اطلاعات و دانش شیمیایی خود، دانستنی‌های بیش‌تری را در مورد سیستم موردنظر مطالعه بدست آورد.

همانطور که می‌دانیم امروزه با توسعه‌ی دستگاه‌های آزمایشگاهی برای هر آزمایش سیل عظیمی از داده بدست می‌آید، این امر سبب شده‌است که بیش‌تر شیمیدانان تجزیه نقش توسعه‌دهنده‌ی روش‌های کمومتریکس و همچنین پیدا کردن کاربردهای جدید آن‌ها در علم شیمی را داشته‌باشند.

نکته‌ای را که پروفیسور سوانته ولد بنیانگذار کمومتریکس و دیگر دانشمندان عرصه‌ی کمومتریکس همواره بدان اشاره داشتند این است که کاربران روش‌های کمومتریکس می‌بایست در خاطر داشته باشند که شیمی به عنوان مرکز و پایه‌ی کمومتریکس است و رهیافت‌های کمومتریکس بدون توجه به اصول شیمی و بدون مورد قبول بودن در حیطه‌ی شیمی، هیچ گونه ارزشی ندارد [۲].

جدا از کاربردهای بهینه‌سازی آزمایش توسط روش‌های کمومتریکس، با توجه به هدف کاربر از روش‌های کمومتریکس، این روش‌ها می‌توانند به شکل‌های توصیفی^۲ و یا پیشگویانه^۳ بکار روند. در شکل توصیفی خواص سیستم‌های شیمیایی به تنها منظور شناخت ساختار نهفته در داده و همچنین رابطه‌ی بین نمونه‌ها، مدل می‌شوند. در شکل پیشگویانه مدل‌سازی به منظور شناخت ساختار، رفتار و خواص نمونه‌های جدید ایجاد می‌شود.

¹ Experimental Design

² Descriptive

³ Predictive

۱.۱.۳ طراحی دارو

با توجه به آنچه که برای تعریف کمومتریکس در بالا اشاره شد، می توان به بالا بودن جایگاه بزرگ کمومتریکس در علم شیمی پی برد. امروزه با تعامل هرچه بیشتر شاخه های مختلف علوم و مخصوصا در بین شیمی و علوم زیستی، کمومتریکس توانسته است که مشکلات زیادی را از پیش روی مطالعات زیستی بردارد. یکی از مواردی که اخیرا به وسیله دانشمندان کمومتریکس توسعه یافته چشمگیری داشته است، مبحث طراحی دارو^۱ و دسته ی QSAR/QSPR (مطالعه ی رابطه ی کمی فعالیت/خاصیت - ساختار) می باشد.

شیمی محاسباتی پتانسیل قابل توجهی برای طراحی دارو دارد، به منظور بهبود شرایط و استفاده ی منطقی از شیمی محاسباتی بهتر است که مدل سازی مولکولی با تحلیل ریاضی رابطه ی بین پاسخ های آشکار شده توسط آزمایش های زیستی^۲ و مجموعه ای مناسب از خواص مولکولی که از روش های محاسباتی بدست آمده اند، همراه باشد. یک راه معمولی قبول شده برای غلبه بر نبود داده های آزمایشگاهی در پدیده های شیمیایی پیچیده نظیر طراحی دارو، آنالیز بر اساس رابطه های کمی بین ساختار و فعالیت (خواص) است.

۱.۲ کمومتریکس در QSAR/QSPR

اکنون با استفاده از تعریف کمومتریکس به چگونگی کاربرد کمومتریکس در QSAR می پردازیم:

- فرض اولیه ی بکار رفته در QSAR بر پایه ی دانش های شیمی و زیستی حدود بیش از یک قرن و نیم پیش در سال ۱۸۴۱ میلادی توسط براون^۳ و فرزر^۴ به این صورت بیان شد که فعالیت زیستی یک ماده به ساختار (ترکیب) آن وابسته است [۵]. در مرحله ی بعد تعدادی مولکول همجنس^۵ انتخاب می شوند، سپس میزان فعالیت زیستی آنها به وسیله ی یک روش معتبر که می تواند با روش طراحی آزمایش بهینه شده باشد، به طور یکسان برای تمام مولکول ها اندازه گیری می شود

¹ Drug Design

² Bioassays

³ A. Crum Brown

⁴ Thomas R. Fraser

⁵ Congeneric

همچنین رمزی کردن^۱ اطلاعات مربوط به ساختار مولکول‌ها که توسط روش‌های ریاضی انجام می‌شود. در مرحله‌ی دوم است که با استفاده از روش‌های گوناگون مدل‌سازی، سعی بر ایجاد یک مدل بهینه و معتبر از نقطه نظر شیمیایی و آماری و با کیفیت از لحاظ پیش‌پیش بینی می‌شود. کیفیت مدل ایجاد شده در مرحله‌ی دوم نه تنها به روش‌های بکاررفته در مدل‌سازی، بلکه به صحت و درستی داده‌های بدست‌آمده از آزمایش در مرحله‌ی اول نیز بستگی دارد. سپس از روی مدل ایجاد شده، سعی در تفسیر رفتار مولکول‌ها بر اساس ساختار آنها می‌شود.

با توجه به موارد گفته‌شده در بالا، نمودار گردش^۲ در شکل (۱-۲) شمایی از خط مشی^۳ QSAR را نمایش می‌دهد.

یکی از محدودیت‌های مدل‌های مورد مطالعه‌ی QSAR، کیفیت داده‌های آزمایشگاهی می‌باشد. کیفیت و صحت مدل ایجاد شده و در نهایت فعالیت پیش‌بینی شده توسط مدل، نمی‌تواند از بیش‌تر از صحت مربوط به داده‌های آزمایشگاهی باشد. همچنین امکان دارد که داده‌های بکاررفته در مدل از منابع مختلف همچون سایر مقالات جمع‌آوری شده‌باشد، بنابراین در این شرایط نیز باید سعی شود برای جلوگیری از گوناگونی و تفاوت در نوع نمونه‌ها، داده‌ها از یک منبع مشخص گرفته شود [۶].

مرحله‌ی نهایی که معمولاً هدف اکثر مطالعات QSAR می‌باشد، انتخاب روشی برای ارزیابی و اعتباربخشی^۴ مدل است که علاوه بر ارزیابی عملکرد برازش^۵ مدل، اجازه‌ی ارزیابی قدرت پیش‌بینی مدل را می‌دهد. معمولاً مورد دوم مهمترین خصوصیت برای یک مدل QSAR قابل قبول مطرح می‌شود.

میزان قدرت پیش‌بینی مدل به وسیله‌ی تقسیم کردن مجموعه‌ی داده‌ها به مجموعه‌ی تمرین که مدل از روی آن محاسبه می‌شود و مجموعه‌ی آزمایش که قدرت پیش‌بینی مدل را ارزیابی می‌کند [۶].

¹ Encoding

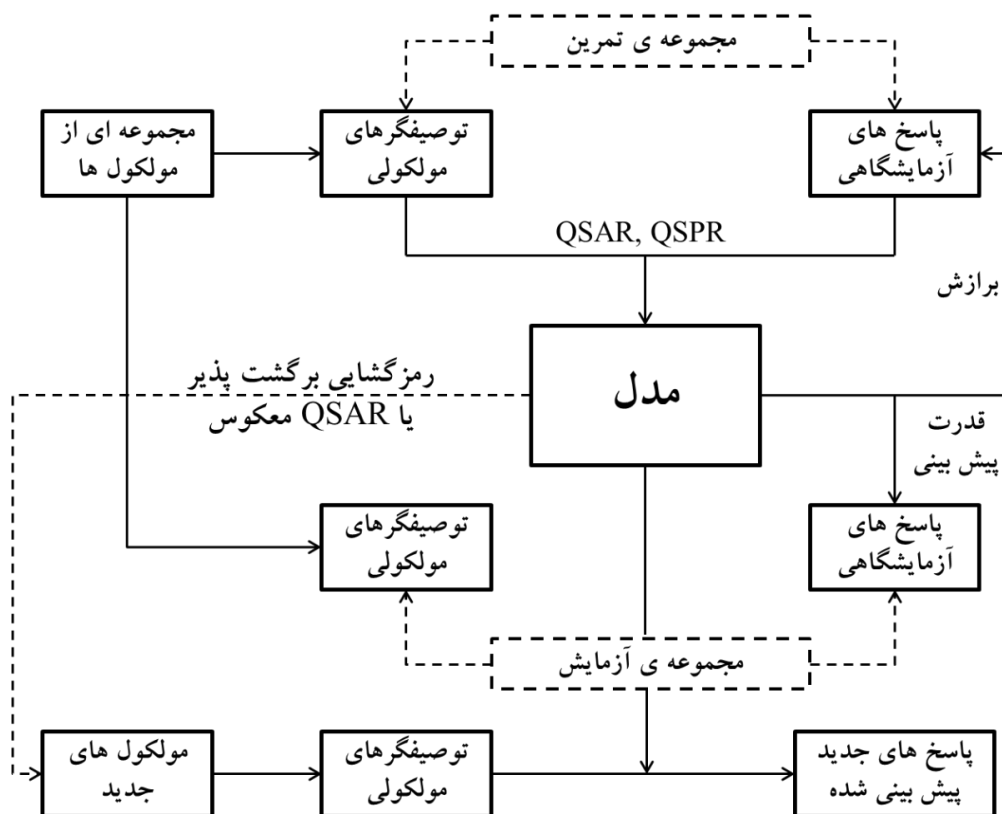
² Flowchart

³ Strategy

⁴ Validation

⁵ Fitting

در حدود نزدیک به ۵۰ سال پیش بود که اولین نمونه‌هایی از کاربرد QSAR (مطالعه‌ی رابطه‌ی کمی ساختار-فعالیت) در مطالعات مربوط به شیمی زراعی، شیمی دارویی، سم‌شناسی و سرانجام در بیش‌تر شاخه‌های شیمی پیداشد [۷].



شکل (۱-۲): شمای عمومی خط مشی QSAR/QSPR [۶]

افزایش چشمگیر کاربرد و توانایی آن را می‌توان از دو جهت مورد بررسی قرار داد [۸]:

- اهمیت فرض اولیه‌ای که نشان می‌دهد فعالیت زیستی مولکول‌ها به صورت تابعی از ساختاری است که به وسیله‌ی ویژگی‌های الکترونیکی، آب‌گریزی و خصوصیات فضایی بیان می‌شود.
- رشد و توسعه‌ی سریع در روش‌ها و ابزارهای محاسباتی که به دنبال تعیین و تصحیح کردن تعداد زیادی از متغیرهایی که مولکول‌ها را توصیف می‌کنند، بوجود آمده است.

با توجه به آنچه که تا به اینجا در مورد QSAR/QSPR توضیح داده شد می توان بیان کرد که این دو رهیافت بر این اصل استوار هستند که ترکیب های متفاوت شیمیایی می توانند به صورت ریاضی گونه ای از خواص مولکولی و یا توصیف گرهای مولکولی رمزگذاری شوند. بنابراین یک روش مدل سازی آماری قادر خواهد بود تا به جستجوی وابستگی موجود در بین توصیف گرهای مولکولی و خصوصیت هدف تعریف شده پرداخته و از آن برای پیش بینی صحیح خصوصیت متناظر ترکیب های جدید استفاده کند. در اینجا اگر هدف نهایی مورد بررسی و پیش بینی (فعالیت مورد نظر مولکول) به صورت مجموعه ای از اعداد پیوسته باشد، قالب رگرسیون برای مدل سازی انتخاب می شود. در حالیکه اگر هدف از رمزگذاری مولکول ها شناسایی دسته های مختلف در ترکیبات باشد از روش های طبقه بندی مجزا^۱ بدین منظور استفاده می شود. (در این مطالعه مدل سازی صرفاً برای پیش بینی فعالیت مولکول های مورد بررسی انجام خواهد گرفت).

در هر رهیافت مبتنی بر QSAR/QSPR، چند مرحله کلیدی با عناوین جمع آوری داده، محاسبه ی توصیف گرهای مولکولی مناسب از روی ساختارهای شیمیایی متناظر، در صورت ضرورت انتخاب آنها، ساخت مدل و در پایان ارزیابی درونی و بیرونی مدل ساخته شده وجود دارد [۹].

در اینجا سعی شده است تا با ترتیب مراحل بیان شده در بالا، داده های مورد استفاده، روش های انتخاب توصیف گر، روش های مختلف مدل سازی برای داده های چند متغیره، روش های ارزیابی مدل و به طور کلی فرایند پیش روی این مطالعه توضیح داده شود. از آنجاییکه در اکثر موارد، روش های انتخاب متغیر با مدل سازی همراه است، در ادامه ابتدا در مورد روش های مدل سازی و سپس روش های انتخاب متغیر بحث خواهد شد.

۱.۳ داده های مورد استفاده در این مطالعه

داده هایی که در این مطالعه از آنها استفاده شده شامل سه مورد گابا^۲، فلور^۳ و ایدز^۴ است که به وسیله ی پژوهشگران تعدادی توصیف گر برای مولکول های مورد نظر محاسبه شده است.

¹ Discrete Classification

² GABA

³ FLUOR

⁴ HIV

FLUOR: شامل ۱۱۶ ترکیب آلی که توسط ۱۲۶۸ توصیفگر نظری بیان شده است. متغیر معیار در این داده میزان فلوئور دوستی این ۱۱۶ ترکیب آلی است [۱۰].

GABA: این داده شامل ۷۸ مشتق از ترکیب فلاون است که به همراه ۱۱۸۷ توصیفگر مولکولی میزان تمایل لیگانهای فلاونوئید را برای قرار گرفتن در محل های بنزودی آزپین پذیرنده های گابا [۱۱] مورد بررسی قرار می دهد [۱۰].

HIV: ۱۰۷ مشتق از داروی HEPT به همراه ۱۶۰ توصیفگر مولکولی این داده را تشکیل می دهند و فعالیت مورد بررسی در این مطالعه نیمه‌ی بیشترین غلظت موثر از این دارو برای حفاظت از سلول های MT-4 در برابر بیماری HIV-1 است [۱۲].

۱. ۴ روش های مدل سازی

داده های اندازه گیری شده همان اطلاعات نیستند. بنابراین یک مورد مهم در تمام علوم تجربی که وظیفه ی اندازه گیری داده ها را به عهده دارد، این است که چگونه اطلاعات مناسب برای هدف آزمایش از داده ها آشکار شود [۱۳]. بدست آوردن اطلاعات دقیق تر که نتیجه ی استفاده ی از داده های چندمتغیره نسبت به تک متغیره است [۱۴] و همچنین توسعه ی روزافزون دستگاه های تجزیه ای باعث شده تا بیان یک نمونه با صدها و یا هزاران متغیر اندازه گیری شده (مانند طول موج، زمان بازداری، نسبت جرم به بار و یا توصیفگرهای شیمیایی) رواج یابد. بنابراین به منظور بدست آوردن اطلاعات مناسب و با دقت کافی در مورد سیستم مورد نظر، می بایست از روش های مدل سازی چند متغیره استفاده کرد. بیشتر مسائل مربوط به مدل سازی در علوم کاربردی می تواند در قالب رگرسیون بیان شود [۱۵]. از این رو یک مدل رگرسیون به منظور ارتباط دادن مجموعه ای از متغیرهای توصیفگر (متغیرهای مستقل)^۱ به یک متغیر معیار (متغیر وابسته)^۲ استفاده می شود.

¹ Independent variables

² Dependent variable

- به طور کلی روش‌های رگرسیون برای رسیدن به دو هدف استفاده می‌شود که در زیر با مثالی بیان شده‌اند:
- کاربرد رگرسیون در فرایندهای صنعتی شیمیایی منحصر به منظور پیش‌بینی است. از آنجاییکه اندازه‌گیری تمامی متغیرهای درگیر در سیستم مورد بررسی در فرایندهای صنعتی آسان نبوده و همراه با پیچیدگی‌های خاصی همراه است، مدل رگرسیون بکار رفته در اینجا به منظور ارتباط دادن متغیرهای اندازه‌گیری شده به متغیری است که تعیین میزان آن اهمیت زیادی در سیستم داشته و اندازه‌گیری مستقیم آن امکان‌پذیر نبوده و یا پیچیده است. بنابراین مدل رگرسیون با ارتباط بین این دو دسته قاعده‌ی پیش‌بینی را بنا می‌نهد و اندازه‌گیری متغیرهای پیچیده را امکان می‌سازد.
 - کاربرد دیگر آن را می‌توان در زمینه‌ی فناوری زیستی^۱ مانند QSAR/QSPR مثال زد، به طوری که در اینجا هدف از مدل‌سازی بر اساس رگرسیون پی بردن به چگونگی تاثیر هر کدام از توصیفگرهای ساختاری بر فعالیت زیستی مولکول مورد نظر و یا تخمین فعالیت مولکول‌های جدید است.

۱.۴.۱ رگرسیون خطی چندتایی^۲:

هدف یک روش رگرسیون ارتباط دادن مجموعه‌ای از متغیرهای مستقل (توصیفگرهای شیمیایی) به متغیر وابسته (متغیر معیار) است که عموماً به شکل خطی زیر بیان می‌شود:

$$y_i = \mathbf{x}_i \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, \dots, n \quad (1-1)$$

در اینجا y_i به عنوان i امین متغیر وابسته (مقدار فعالیت و یا خاصیت)، \mathbf{x}_i ، i امین سطر از ماتریس داده که رمزی شده‌ی یک مولکول به شکل عددی است. $\boldsymbol{\beta}$ بردار ضرایب رگرسیون و هم بعد با بردار \mathbf{x}_i و ε_i مقدار خطای مدل در پیش‌بینی مربوط به فعالیت i امین مولکول است.

همانطور که می‌دانیم بردار $\boldsymbol{\beta}$ در جهتی بدست می‌آید که مربع تفاوت میزان پیش‌بینی شده‌ی برای بردار \mathbf{y} توسط مدل و میزان حقیقی آن کمترین مقدار شود که توسط تابع فقدان^۳ کمترین مربعات بیان می‌شود:

¹ Biotechnology

² Multiple Linear Regression

³ Loss function

$$f(\beta) = \|y - X\beta\|^2 \quad (2-1)$$

اگر ماتریس $X'X$ غیرمنفرد^۱ باشد، بردار β به صورت منحصر بفرد^۲ و با استفاده از رابطه‌ی زیر بدست می‌آید.

$$\beta = (X'X)^{-1}X'y \quad (3-1)$$

اما اگر در شرایطی ماتریس $X'X$ منفرد^۳ باشد (وجود وابستگی در بین متغیرها و ترکیب خطی بودن متغیرها از یکدیگر)، مجموعه‌ای از بردارهای β بدست خواهد آمد که تمامی آنها برای معادله‌ی (۲-۱) جواب هستند و از رابطه‌ی زیر بدست می‌آیند:

$$\beta = X^+y + NV \quad (4-1)$$

که در این رابطه بالانشان + نماد معکوس مور- پنروز^۴ است، N ماتریس متمم عمود^۵ X' است به طوریکه $XN = 0$ و V ماتریسی دلخواه با اندازه‌ی مناسب است [۱۵]. بنابراین با انتخاب‌های متفاوت از V ، بردارهای متفاوتی از β بدست خواهد آمد که تمامی آنها معیار معادله‌ی (۲-۱) را دارا هستند. در این مورد یک راه معمول برای انتخاب از بین این مجموعه‌ی بردارها این است که مقدار $V = 0$ را انتخاب کنیم.

در واقع می‌توان رابطه‌ی (۳-۱) را نیز به شکل X^+y بیان نمود، زیرا هنگامی که $X'X$ غیرمنفرد باشد آنگاه $X^+ = (X'X)^{-1}X'$ است. هنگامی که ماتریس $X'X$ منفرد بدست می‌آید و یا به عبارتی بعضی از مقادیر تکین^۶ آن بسیار کوچک و نزدیک صفر است، رویه‌ی معکوس مور- پنروز سبب می‌شود که معکوس آن مقادیر تکین و به تبع از آن عناصر β بسیار بزرگ شوند و جواب‌های بسیاری نیز برای β وجود داشته‌باشد. به طور معمول دیده شده‌است که بزرگ شدن عناصر بردار β باعث کاهش عمومیت این بردار برای استفاده‌ی در مجموعه-های دیگری از داده می‌شود [۱۵].

¹ Nonsingular

² Unique

³ Singular

⁴ Moore-Penrose Inverse

⁵ Orthocomplement (Orthogonal Complement)

⁶ Singular Values

۱. ۴. ۲ رگرسیون مرزی^۱:

یک روش معمول که برای جلوگیری از بزرگ شدن اندازه‌ی بردار β استفاده از انقباض^۲ است که به صورت زیر در هنگامی که ماتریس $X'X$ منفرد بدست می‌آید، بکار می‌رود:

$$g(\beta) = \|y - X\beta\|^2 + \lambda\|\beta\|^2 \quad (5-1)$$

فرایند بکار رفته در رابطه‌ی (۵-۱) به رگرسیون مرزی معروف است. عدد λ بیان کننده‌ی میزان جریمه^۳ برای کاهش اندازه‌ی بردار β است. کمترین مقدار برای رابطه‌ی (۵-۱) از رابطه‌ی زیر بدست می‌آید:

$$\beta = (X'X + \lambda I)^{-1} X'y \quad (6-1)$$

توجه داشته باشید که تمامی مقادیر منفرد ماتریس معکوس در این رابطه همواره غیر صفر خواهند بود.

با این حال انقباض صورت گرفته بر روی بردار β با روش رگرسیون مرزی جوابگوی بهترین جواب برای مدل‌سازی بر اساس رگرسیون نیست و نمی‌تواند پشتوانه‌ای برای مجموعه‌ی دیگری از داده باشد. رگرسیون مرزی نوع دیگر از یک روش کلی‌تر به نام تنظیم تیخونوف^۴ است [۱۶]. رهیافت دیگری که به منظور بهینه‌کردن فرایند رگرسیون به نظر می‌رسد، استفاده از خلاصه‌ی متغیرهای موجود در ماتریس X است.

۱. ۴. ۳ روش‌های رگرسیون چندمتغیره مبتنی بر متغیرهای پنهان

معمول‌ترین روش برای بدست آوردن خلاصه‌ای از متغیرها، استفاده از روش تحلیل اجزای اصلی^۵ (PCA) است. ماتریس X با استفاده از PCA به دو زیرماتریس امتیاز^۶ (T) و ظرفیت^۷ (P) تجزیه می‌شود و هر کدام شامل اجزای اصلی ماتریس X هستند که به صورت ترکیب خطی ستون‌ها در زیرماتریس امتیاز و ترکیب خطی سطرها

¹ Ridge Regression

² Shrinkage

³ Penalty

⁴ Tikhonov regularization

⁵ Principal Components Analysis

⁶ Score

⁷ Loading

در وزن آورده شده است و به اصطلاح به آنها متغیرهای پنهان^۱ نیز می‌گویند. حال که متغیرهای موجود در ماتریس X توسط روش PCA به تعدادی محدود از متغیرهای پنهان خلاصه شد، می‌توان از این متغیرها به عنوان متغیرهای جدیدی که متغیر معیار را توصیف می‌کنند در یک مدل رگرسیون استفاده کرد که روش‌های رگرسیون اجزای اصلی^۲ (PCR) و کمترین مربعات جزئی^۳ (PLS) بر مبنای این متغیرها عمل می‌کنند.

هر دوی این روش‌ها با استفاده از مختصات جدیدی که توسط متغیرهای پنهان برای سطر و ستون ماتریس X توصیف می‌کنند، سعی در یافتن بهترین جواب برای β هستند. با توجه به محدود بودن تعداد متغیرهای پنهان (ماتریس T) و عمود بودن آنها، ماتریس $T'T$ غیرمنفرد بدست می‌آید و تنها یک جواب بهینه برای β وجود خواهد داشت. تفاوت این دو روش در چگونگی بدست آوردن ضرایب رگرسیون و یا به عبارتی بدست آوردن متغیرهای پنهان (ماتریس T) است. در روش PCR این متغیرها با استفاده از روش PCA بدست می‌آیند. یکی از مواردی که علاوه بر ورود متغیرهای پنهان به مدل در مقایسه با متغیرهای اصلی به بهبود مدل ساخته شده کمک می‌کند، خلاصه شدن ماتریس داده در تعداد محدودی متغیر پنهان است که از ورود خطای مربوط به داده‌های اندازه‌گیری شده در مدل رگرسیون تا حدود زیادی جلوگیری می‌کند.

در روش PLS علاوه بر وجود تعداد کمتر متغیرهای پنهان بکاررفته در مدل و بهینه بودن ضرایب رگرسیون، دو ویژگی دیگر وجود دارد که باعث برتری آن به روش PCR می‌شود. از آنجاییکه در روش PCR همانطور که در بالا اشاره شد تا حدودی مدل ساخته شده از خطای همراه با داده‌ی اندازه‌گیری شده در امان است، در روش PLS به این علت که برای هر دو مجموعه‌ی متغیرهای مستقل و وابسته مجموعه‌ی متغیرهای پنهان بدست می‌آید، کاهش خطا را در هر دو ماتریس خواهیم داشت. در ضمن در روش PCR از آنجاییکه هیچ‌گونه وابستگی بین خطای همراه با ماتریس داده‌های مستقل و وابسته وجود ندارد، متغیرهای پنهان بدست آمده برای ماتریس مستقل در جهتی تصادفی و غیر از جهت متغیرهای پنهان ماتریس وابسته بدست می‌آیند. این تفاوت در جهت منجر به زاویه‌دار شدن دو بردار نسبت بهم می‌شود و مدل رگرسیون قابل قبولی را ایجاد نمی‌کند. در روش PLS متغیرهای پنهان در دو ماتریس داده‌ی وابسته و مستقل طوری در فضای پنهان چرخیده می‌شوند که این زاویه به

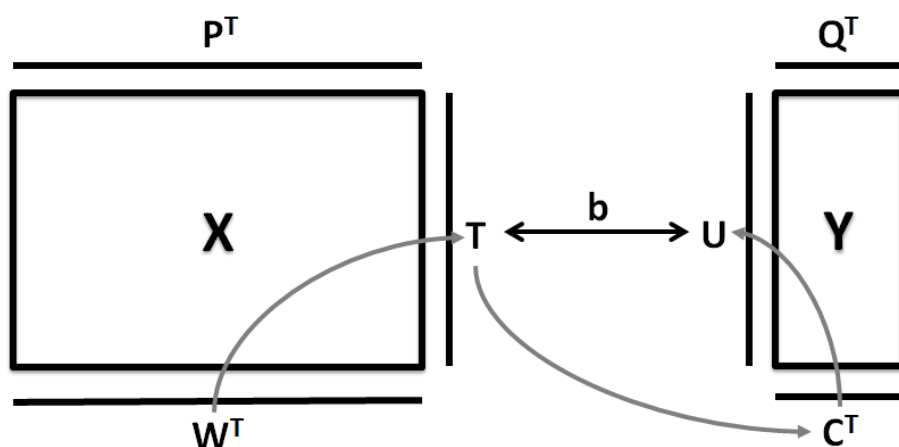
¹ Latent Variables

² Principle Components Regression

³ Partial least Squares

کمترین مقدار خود برسد و به عبارتی فرایند حذف خطا به مراتب خیلی بهتر از روش **PCR** انجام می‌شود. شمایی از چگونگی رویه‌ی روش **PLS** در شکل (۳-۱) آورده شده‌است.

بردارهای **t** و **u** متغیرهای پنهانی هستند که از رویه‌ی **PLS** به ترتیب برای ماتریس‌های **X** و **y** بدست می‌آیند، اطلاعات تمام ستون‌های **X** و **y** را دارا هستند و همچنین مقداری از پراکنش^۱ غیرمشترک در بین **X** و **y** را حذف کرده‌اند. به ازای تعداد متغیر پنهان در نظر گرفته برای مدل **PLS** از بردارهای **t** و **u** بدست می‌آید و در نهایت یک رگرسیون خطی چندمتغیره بین ماتریس‌های **T** و **U** به منظور مدل‌سازی اطلاعات بین **X** و **y** انجام می‌شود. همانطور که پیش‌تر بدان اشاره شد، متغیرهای پنهان بکاررفته در روش **PLS** طوری چرخیده می‌شوند که بیشترین همبستگی بین **T** و **U** بدست‌آید. این چرخیدن به این صورت انجام می‌شود که ماتریس **T** از ترکیب خطی ستون‌های **X** و **w** بدست می‌آید که خود بردار وزن **w** از ترکیب خطی ستون‌های ماتریس **X** و بردار **y** بدست می‌آید.



شکل (۳-۱): شمایی کلی رابطه‌ی بین ماتریس‌های **X** و **y** در رویه‌ی **PLS**

اگرچه روش‌های بر پایه‌ی متغیرهای پنهان در ابتدا در روانشناسی معرفی شدند، امروزه در بسیاری از زمینه‌ها مانند علوم اجتماعی، علوم طبیعی و اقتصاد استفاده می‌شوند. ورود این روش‌ها را از روانشناسی به دیگر رشته‌ها را می‌توان از اواخر دهه‌ی پنجاه میلادی تخمین زد.

¹ Variance