

دُنْيَا



پایان نامه دوره کارشناسی ارشد در رشته شیمی گرایش تجزیه

عنوان:

**مدل سازی QSPR پتانسیل اکسایش-کاهش
آنتری اکسیدان های فنولی**

استاد راهنما:

دکتر محمد حسین فاطمی

استاد مشاور:

دکتر محمدرضا حاج محمدی

اساتید داور:

دکتر محمد جواد چایچی دکتر سید ناصر عزیزی

دانشجو:

زهره قره چاهی

شهریور ماه ۱۳۹۰

پاسکناری

ستایش و پاس بیکران بایست آن ایزد دانایی است که چراغ داش را دانیده انسان فروزان می دارد، تا در پرتو آن خود را از برگ کی جل برقند و به آزادگی و

برفروزی دست یابد.

از پر و ماد عزیزم

که کوهر و جود شان، نیم کلاشان، باران محبت شان و دستان همیشه یار یکر شان هموار کننده راه زندگی ام بوده است، پاسکنارم و بر دستان پر مهر شان بوسه
می زنم.

از برادر مهر بام، احسان، که در دوران تحصیل همواره راهنمای پشتیبانم بوده، بی نهایت پاسکنارم و بهترین هزار ایش آرزومندم.

در اینجا وظیفه شاگردی خود می دانم تا از زحات استاد بزرگوارم جناب آقای دکتر محمد حسین فاطمی، که مرآبه عنوان عضوی از گروه خود پذیر فتنه و دظام
مراحل انجام پیمان نامه صبورانه راهنمای پشتیبان من بودند، نهایت قدر رانی را داشته باشم. یعنی می سیماز ترین شکرات خود را به استاد مشاور محترم جناب
آقای دکتر محمد رضا حاج محمدی تقدیم می دارم. از استادید او را جناب آقای دکتر محمد جواد چایچی و جناب آقای دکتر سید ناصر عزیزی که زحمت مطالعه پژوهه
بنده را تقبل نمودند بسیار مشکرم.

دیگران از دوست بسیار مهر بام، یا نهایت آشنایی با وی دنیایی جدیدی را پیش روی چشم انداختم که شود و داین دوره از تحصیل همواره مشوق و حامی من بوده،

بی نهایت پاسکنارم.

تعدیم به

عزیزترین مهیت‌های زندگیم

مروارید
•

و دردانه‌های همیشه قلبم

احسان، ناہید

ایمان و نادیا

چکیده

روابط کمی ساختار-فعالیت یا ویژگی (QSAR/QSPR) یکی از فنون نویدبخش در زمینه روش‌های مجازی به منظور پیش‌بینی ویژگی‌های شیمیایی است. این روش‌ها، با استفاده از توصیف کننده‌هایی که از ساختار مولکولی منتج می‌شوند، به جستجوی الگویی در داده‌ها می‌پردازند تا فعالیت یا ویژگی مواد شیمیایی جدیدی را که ویژگی‌های مولکولی مشابهی دارند، پیش‌بینی کنند. در بخش اول این پژوهه، از روش QSPR، جهت پیش‌بینی پتانسیل اکسایش-کاهش ۴۲ آنتی‌اکسیدان فولی استفاده شد. به منظور انتخاب مهم‌ترین توصیف کننده‌ها از روش رگرسیون خطی چندگانه مرحله‌ای استفاده گردید. سپس، روش رگرسیون خطی چندگانه (MLR) و شبکه عصبی پرسپترون چند لایه (MLP NN) جهت ساخت مدل‌های QSPR خطی و غیرخطی به کار گرفته شد. مقایسه نتایج آماری این دو مدل نشان داد که مدل MLP NN در پیش‌بینی پتانسیل اکسایش-کاهش آنتی‌اکسیدان‌های فولی از اعتبار بیشتری برخوردار است. بعلاوه، بررسی توصیف کننده‌های موجود در مدل‌های QSPR نشان داد که خصوصیات الکترونی مشتقات فولی، نقش مهمی در ویژگی‌های آنتی‌اکسیدانی آن‌ها دارد.

در بخش دوم این پژوهه، مدل‌های QSPR بر پایه روش‌های MLR، MLP NN و رگرسیون بردار پشتیبان (SVR)، جهت پیش‌بینی نیمه‌عمر پالایشی ۶۲ مولکول با فنیل چندکلره در ماهی قزل‌آلای رنگین‌کمان ساخته شد. در این بخش از پژوهه، از الگوریتم رنگیک به عنوان روش انتخاب متغیر استفاده شد. مقادیر عددی بزرگ Q_{LOO}^2 و R^2 و مقدار کم RMSE، برتری مدل MLP NN و همچنین وابستگی غیرخطی ویژگی‌های ساختاری مولکولی به نیمه‌عمر پالایشی با فنیل‌های چندکلره را تأیید می‌کند. تجزیه و تحلیل توصیف کننده‌های موجود در مدل‌ها حاکی از این است که ویژگی‌های ساختاری ۲ بعدی مولکول، تراکم و الکترونگاتیویته از عوامل اصلی در تعیین نیمه‌عمر پالایشی ترکیبات با فنیل‌های چندکلره می‌باشد.

واژه‌های کلیدی: پتانسیل اکسایش-کاهش، آنتی‌اکسیدان‌های فولی، نیمه‌عمر پالایشی، با فنیل‌های چندکلره، رابطه کمی ساختار ویژگی، شبکه عصبی مصنوعی

فهرست مطالب

صفحه	عنوان
۱	فصل اول: مقدمه
۳	۱-۱- روش‌های QSAR/QSPR
۴	۲-۱- اهداف روش‌های QSAR/QSPR
۵	۳-۱- روش‌های مدل‌سازی
۷	فصل دوم: تئوری
۹	۱-۲- جمع‌آوری و انتخاب سری داده‌ها
۱۰	۲-۲- رسم و بهینه‌سازی ساختارهای مولکولی
۱۱	۳-۲- محاسبه توصیف کننده‌های مولکولی
۱۱	۱-۳-۲- توصیف کننده‌های تجربی
۱۲	۲-۳-۲- توصیف کننده‌های تئوری
۱۲	۱-۲-۳-۲- توصیف کننده‌های ساختاری
۱۳	۲-۲-۳-۲- توصیف کننده‌های توپولوژیکی
۱۳	۱-۲-۲-۳-۲- توصیف کننده‌های جزء
۱۳	۲-۲-۲-۳-۲- شاخص‌های توپولوژی
۱۴	۳-۲-۲-۳-۲- توصیف کننده‌های زیرساختاری
۱۴	۴-۲-۲-۳-۲- توصیف کننده‌های محیطی
۱۵	۳-۲-۳-۲- توصیف کننده‌های هندسی
۱۶	۴-۲-۳-۲- توصیف کننده‌های الکترونی
۱۷	۴-۲- تجزیه و تحلیل آماری توصیف کننده‌ها و انتخاب مؤثرترین آنها
۱۸	۱-۴-۲- رگرسیون خطی چندگانه مرحله‌ای
۱۹	۲-۴-۲- الگوریتم ژنتیک
۲۱	۱-۲-۴-۲- مفاهیم و اصول الگوریتم ژنتیک
۲۶	۲-۵- ایجاد مدل‌های آماری
۲۷	۱-۵-۲- رگرسیون خطی چندگانه
۲۸	۲-۵-۲- شبکه‌های عصبی مصنوعی
۲۸	۱-۲-۵-۲- ساختار سلول عصبی انسان
۳۰	۲-۲-۵-۲- عنصر پردازش در شبکه‌های عصبی مصنوعی
۳۳	۳-۲-۵-۲- ساختار شبکه‌های عصبی
۳۶	۴-۲-۵-۲- انواع شبکه‌های عصبی

۳۶ آموزش شبکه‌های عصبی ۲-۵-۲
۳۹ ۲-۵-۳- ماشین بردار پشتیبان.....
۴۰ ۲-۵-۳-۱- اصول رگرسیون بردار پشتیبان.....
۴۴ ۲-۵-۳-۲- انواع توابع کرnel.....
۴۵ ۲-۵-۳-۳- آموزش رگرسیون بردار پشتیبان.....
۴۷ ۲-۶- تحلیل و ارزیابی آماری مدل‌ها و انتخاب بهترین مدل.....
۵۱ ۲-۷- نرم‌افزارهای مورد استفاده.....
۵۱ ۲-۷-۱- بسته نرم‌افزاری HyperChem.....
۵۲ ۲-۷-۲- بسته نرم‌افزاری DRAGON.....
۵۲ ۲-۷-۳- بسته نرم‌افزاری MOPAC.....
۵۳ ۲-۷-۴- بسته نرم‌افزاری CODESSA.....
۵۳ ۲-۷-۵- بسته نرم‌افزاری SPSS.....
۵۴ ۲-۷-۶- بسته نرم‌افزاری STATISTICA.....
۵۴ ۲-۷-۷- بسته نرم‌افزاری MATLAB.....

۵۵ فصل سوم: مدل‌سازی QSPR پتانسیل اکسایش-کاهش آنتی‌اکسیدان‌های فنولی.....
۵۹ ۳-۱- روش کار.....
۵۹ ۳-۱-۱- سری داده‌ها.....
۶۱ ۳-۱-۲- محاسبه و پیش‌پردازش توصیف کننده‌ها.....
۶۲ ۳-۱-۳- انتخاب توصیف کننده‌ها و مدل‌سازی خطی.....
۶۶ ۳-۱-۴- مدل‌سازی غیرخطی با شبکه‌ی عصبی پرسپترون چند لایه.....
۶۷ ۳-۲- بحث و نتیجه‌گیری.....
۶۷ ۳-۲-۱- بررسی نتایج.....
۷۲ ۳-۲-۲- بررسی و تفسیر توصیف کننده‌ها.....
۷۳ ۳-۳- نتیجه‌گیری کلی.....

۷۴ فصل چهارم: مدل‌سازی و پیش‌بینی نیمه‌عمر آلاینده‌های با فنیل‌های چند کلره با بهره‌گیری از رویکرد QSPR.....
۷۸ ۴-۱- روش کار.....
۷۸ ۴-۱-۱- سری داده‌ها.....
۸۰ ۴-۱-۲- محاسبه و پیش‌پردازش توصیف کننده‌ها.....
۸۱ ۴-۱-۳- انتخاب توصیف کننده‌ها و مدل‌سازی خطی.....
۸۳ ۴-۱-۴- مدل‌سازی غیرخطی.....
۸۳ ۴-۱-۴-۱- مدل‌سازی غیرخطی با شبکه‌ی عصبی پرسپترون چند لایه.....

۸۵	-۱-۴-۲-۴-۴-۱-۴ مدل سازی غیر خطی با رگرسیون بردار پشتیبان
۸۶	-۴-۲-۴ بحث و نتیجه گیری
۸۶	-۴-۲-۱-۱-۲-۴ بررسی نتایج
۸۹	-۴-۲-۲-۲-۴ ارزیابی داخلی مدل ها
۹۰	-۴-۲-۳-۳-۲-۴ بررسی و تفسیر توصیف کننده ها
۹۲	-۴-۳-۳-۲-۴ نتیجه گیری کلی
۹۳	پیشنهادات برای پژوهش های آینده

فهرست شکل‌ها

صفحه	عنوان
۲۱	شکل ۱-۲- نمایی از نقاط بهینه محلی و کلی
۲۴	شکل ۲- طرحواره‌ای از ترکیب تک نقطه در الگوریتم ژنتیک
۲۵	شکل ۲- طرحواره‌ای از عملگر جهش در الگوریتم ژنتیک
۲۶	شکل ۲- نمودار گردشی الگوریتم ژنتیک
۲۹	شکل ۲-۵- اجزای تشکیل دهنده یک نرون زیستی
۳۰	شکل ۲-۶- شبیه‌سازی نرون مصنوعی از روی نرون زیستی
۳۱	شکل ۲-۷- طرحواره‌ای از عنصر پردازش در شبکه‌های عصبی
۳۳	شکل ۲-۸- تابع فعال‌سازی سیگموئیدی (S-شکل)
۳۴	شکل ۲-۹- طرحواره‌ای از ساختار شبکه عصبی سه لایه
۳۵	شکل ۲-۱۰- ساختار شبکه‌های (الف) پیشخور (ب) برگشتی
۳۸	شکل ۲-۱۱- طرحواره‌ای از ساختار آموزش با ناظر
۴۲	شکل ۲-۱۲- تابع اتلاف و موقعیت بردارهای پشتیبان بر روی آن
۴۴	شکل ۲-۱۳- نگاشت داده‌ها به فضایی با ابعاد بالاتر (0) توسط تابع کرنل
۴۶	شکل ۲-۱۴- تأثیر افزایش ۶ بر روی تعداد بردارهای پشتیبان
۶۷	شکل ۳-۱- طرحواره‌ای از شبکه عصبی پرسپترون سه لایه بهینه شده
۶۸	شکل ۳-۲- نمودار میله‌ای پارامترهای آماری مدل MLP و NN
۶۹	شکل ۳-۳- نمودار مقادیر پیش‌بینی شده پارامتر E_7 با استفاده از مدل MLP NN بر حسب مقادیر تجربی
۷۰	شکل ۳-۴- نمودار باقیمانده مقادیر پارامتر E_7 پیش‌بینی شده با استفاده از مدل MLP NN بر حسب مقادیر تجربی
۷۱	شکل ۳-۵- نتایج تحلیل حساسیت
۸۴	شکل ۴-۱- طرحواره‌ای از شبکه عصبی پرسپترون سه لایه بهینه شده
۸۷	شکل ۴-۲- نمودار میله‌ای پارامترهای آماری مدل MLP و NN
۸۸	شکل ۴-۳- نمودار مقادیر پیش‌بینی شده لگاریتم نیمه‌عمر با استفاده از مدل MLP NN بر حسب مقادیر تجربی
۸۸	شکل ۴-۴- نمودار مقادیر باقیمانده لگاریتم نیمه‌عمر با استفاده از مدل MLP NN بر حسب مقادیر تجربی
۸۹	شکل ۴-۵- نتایج تحلیل حساسیت

فهرست جدول‌ها

عنوان	صفحه
جدول ۱-۳- ترکیبات سری داده و مقادیر تجربی پتانسیل اکسایش-کاهش آن‌ها بر حسب ولت (E_7)	۶۰
جدول ۲-۳- مشخصات توصیف کننده‌های انتخاب شده جهت مدل‌سازی E_7	۶۳
جدول ۳-۳- ماتریس ضرایب همبستگی بین توصیف کننده‌های انتخاب شده	۶۳
جدول ۴-۳- مقادیر عددی توصیف کننده‌ها جهت مدل‌سازی پارامتر E_7	۶۴
جدول ۵-۳- مقادیر تجربی و پیش‌بینی شده E_7 با استفاده از دو مدل MLP و NN با همراه مقادیر باقیمانده آن‌ها	۶۵
جدول ۶-۳- پارامترهای آماری حاصل از مدل‌های MLR و NN	۶۸
جدول ۷-۳- نتایج آزمون درهم آمیختگی- Y	۷۰
جدول ۱-۴- مقادیر تجربی و پیش‌بینی شده نیمه عمر بالایشی ترکیبات با فیل‌های چند کلره	۷۹
جدول ۲-۴- نام و علامت اختصاری توصیف کننده‌های انتخاب شده	۸۲
جدول ۳-۴- ماتریس ضرایب همبستگی بین توصیف کننده‌های انتخاب شده	۸۲
جدول ۴-۴- پارامترهای آماری حاصل از مدل‌های MLR، NN و SVR	۸۶

لیست عالیم و اختصارات

QSAR	رابطه کمی ساختار-فعالیت (Quantitative Structure-Activity Relationship)
QSPR	رابطه کمی ساختار- свойگی (Quantitative Structure-Property Relationship)
MLR	رگرسیون خطی چندگانه (Multiple Linear Regression)
ANN	شبکه عصبی مصنوعی (Artificial Neural Network)
SVM	ماشین بردار پشتیبان (Support Vector Machine)
MLP NN	شبکه عصبی پرسپترون چند لایه (Multi-Layer Perceptron Neural Network)
SVR	رگرسیون بردار پشتیبان (Support Vector Regression)
EVA	توصیف کننده‌های تحلیل مقدار ویژه (Eigen Value Analysis descriptors)
WHIM	توصیف کننده‌های مولکولی کلی نامتغیر وزن دار شده (Weighted Holistic Invariant Molecular descriptors)
3D-MoRSE	توصیف کننده‌های ترکیب هندسه، توپولوژی و وزن‌های اتمی (3D-Molecular Representation of Structure based on Electron diffraction descriptors)
GETAWAY	(GEometry, Topology and Atom-Weights Assembly descriptors)
HOMO	بالاترین اوربیتال مولکولی اشغال شده (Highest Occupied Molecular Orbital)
nOH	تعداد گروه‌های هیدروکسیل (number of hydroxyl group)
LUMO	پایین‌ترین اوربیتال مولکولی اشغال نشده (Lowest Unoccupied Molecular Orbital)
KKT	کاروش-کان-تاکر (Karush-Kuhn-Tucker)
RSS	مجموع مربعات باقیمانده‌ها (Residual Sum of Squares)
TSS	مجموع مربعات کل (Total Sum of Squares)
RMSE	ریشه میانگین مربع خطأ (Root Mean Square Error)
AAE	میانگین قدر مطلق خطأ (Average Absolute Error)
HOMA	مدل نوسانگر هماهنگ آروماتیسیته (Harmonic Oscillator Model of Aromaticity)
IP	پتانسیل یونیزاسیون (Ionization Potential)
PCBs	بای‌فیل‌های چند‌کلره (PolyChlorinated Biphenyls)
BEHm3	بالاترین مقدار ویژه n.3 از ماتریس بردن / وزن دار شده با جرم‌های اتمی (Highest eigenvalue n.3 of Burden matrix/weighted by atomic masses)
Lop	شاخص مرکزی جدا شده (Lopping centric index)
JGI1	میانگین شاخص بار توپولوژیکی از مرتبه ۱ (Mean topological charge index of order 1)
GATS8	خودهمبستگی گری با فاصله ۸ / وزن دار شده با الکترونگاتیویته اتمی ساندرسون (Geary autocorrelation lag-8/weighted by atomic Sanderson electronegativities)

فصل اول

مقدمہ

واژه کمومتریکس^۱ یا شیمی‌سنجدی برای اولین بار در سال ۱۹۷۲ توسط شیمیدانی سوئدی به نام ولد^۲ پیشنهاد گردید و برای نخستین بار در سال ۱۹۷۵ در مجله "علوم کامپیوتر و اطلاعات شیمیایی" توسط کوالسکی^۳ معرفی شد [۱,۲]. کوالسکی تعریف جامعی از علم شیمی‌سنجدی به صورت زیر ارائه داد [۳]:

"شیمی‌سنجدی شاخه‌ای از علم شیمی است که از علوم ریاضی، آمار و منطق جهت طراحی و انتخاب روش‌های بهینه آزمایشگاهی، تحلیل داده‌های شیمیایی به منظور استخراج حداقل اطلاعات شیمیایی و دستیابی به اطلاعات در مورد سیستم‌های شیمیایی استفاده می‌کند."

اساساً هدف شیمی‌سنجدی بهبود فرآیندهای اندازه‌گیری و استخراج اطلاعات شیمیایی مفید از داده‌های اندازه‌گیری شده فیزیکی و شیمیایی می‌باشد. اگر چه شانه شیمی‌سنجدی در علم شیمی تجزیه ایجاد شد اما امروزه قلمرو کاربرد آن بسیار گسترده شده است [۴]. بسیاری از شاخه‌های شیمی مانند شیمی آلی، شیمی معدنی، شیمی محیط‌زیست، شیمی غذایی، شیمی کشاورزی و حتی سایر علوم مانند طراحی دارو،

^۱ Chemometrics

^۲ Wold

^۳ Kowalski

زیست‌شناسی (ژنومیکس^۱، پروتئومیکس^۲ و متابولومیکس^۳) و سم‌شناسی جهت تفسیر داده‌ها به علم شیمی‌سنگی نیازمندند. روش‌های شناخت الگوهای^۴، پردازش علائم^۵، دسته‌بندی^۶، بهینه‌سازی^۷، رابطه کمی ساختار-فعالیت^۸ (QSAR) و رابطه کمی ساختار-ویژگی^۹ (QSPR)، برخی از کاربردهای مهم شیمی‌سنگی در علوم مختلف می‌باشد [۴,۵].

۱-۱- روش‌های QSAR/QSPR

در روش‌های QSAR/QSPR به کمک یک مدل محاسباتی، رابطه‌ای ریاضی بین خصوصیات ساختاری ترکیبات و یک ویژگی اندازه‌گیری شده آن‌ها برقرار می‌گردد [۶]. به منظور ایجاد این رابطه کمی، ویژگی‌های فیزیکو‌شیمیایی (از قبیل چربی‌دوستی، آبگریزی، حلالیت، قطبش‌پذیری) و ویژگی‌های فضایی، الکترونیکی و ساختاری ترکیبات با استفاده از الگوریتم‌های مختلفی به مقادیر عددی تبدیل می‌شود. این کمیت‌های عددی، توصیف‌کننده‌های مولکولی^{۱۰} نامیده شده و به عنوان متغیرهای مستقل در مدل QSAR/QSPR مورد استفاده قرار می‌گیرند. بدین ترتیب، با بهره‌گیری از داده‌های تجربی گروهی از

^۱ Genomics

^۲ Proteomics

^۳ Metabolomics

^۴ Pattern recognition

^۵ Signal processing

^۶ Classification

^۷ Optimization

^۸ Quantitative Structure-Activity Relationship

^۹ Quantitative Structure-Property Relationship

^{۱۰} Molecular descriptors

ترکیبات و توصیف کننده‌های مولکولی آنها، مدلی جهت پیش‌بینی رفتار مولکول‌های مشابه و جدید که خصوصیات ساختاری مشترکی دارند، ساخته می‌شود.

اولین مشاهدات در مورد روابط QSAR/QSPR در رشتۀ سم‌شناسی توسط کراس^۱ صورت گرفت [۷]. وی در سال ۱۸۶۳ در پایان‌نامه دوره دکتری خود، گزارشی در خصوص وجود رابطه بین خصوصیات زیستی و مولکولی ارائه داد. کراس مشاهده کرد که سمیت الكل‌های آلیفاتیک نوع اول با کاهش حلالیت آنها در آب افزایش می‌یابد. به دنبال آن محققین دیگری نیز به وجود روابط خطی بین چربی‌دوستی و ویژگی‌های زیستی مثل سمیت پی بردند [۷]. در این میان، هانش^۲ را می‌توان پیشگام و پایه‌گذار روش QSAR مدرن و به تبع آن QSPR دانست [۸]. وی در اوایل دهه ۱۹۶۰ استفاده از روش رگرسیون خطی چندگانه^۳ (MLR) را به منظور پیش‌بینی فعالیت زیستی ترکیباتی که تا آن زمان سنتز نشده بودند ارائه کرد.

۱-۲- اهداف روش‌های QSAR/QSPR

همان‌گونه که گفته شد روش‌های QSAR/QSPR کاربرد گسترده‌ای در بسیاری از زمینه‌های علمی دارند. در این روش‌ها به طور عمدۀ دو هدف اصلی دنبال می‌شود [۶]:

۱- پیش‌بینی فعالیت‌های زیستی و ویژگی‌های فیزیکو‌شیمیایی مولکول‌ها

۲- درک و توجیه چگونگی عملکرد و مکانیسم مولکول‌ها در فرآیندهای مختلف

در کنار این اهداف، دلایل دیگری سبب توسعه این روش‌ها شده است [۶]:

^۱ Cros

^۲ Hansch

^۳ Multiple Linear Regression

۱- صرفه جویی در وقت و هزینه به منظور تولید و توسعه محصولات (به طور مثال فراورده‌های دارویی و

حشره‌کش‌ها)

۲- کاهش نیاز به آزمایش‌های هزینه‌بر و طولانی مدت بر روی حیوانات و در نتیجه کاهش استفاده از

حیوانات

۳- دستیابی به اهداف شیمی سبز به منظور افزایش کارآیی و یا حذف مواد زائد

۱-۳- روش‌های مدل‌سازی

به منظور ساخت مدل‌های QSAR/QSPR از روش‌های آماری و ریاضی مختلفی استفاده می‌شود.

روش MLR، یکی از انواع روش‌های متداول آماری تحلیل چند متغیره در مطالعات QSAR/QSPR است

[۹]. در این روش تنها روابط خطی پدیده‌های مورد مطالعه بررسی و مدل‌سازی می‌شود. این در حالی است

که بسیاری از ویژگی‌های مولکول‌ها از روابط غیرخطی پیروی می‌کنند. پیشرفتهای اخیر در زمینه علم آمار

منجر به توسعه الگوریتم‌های بسیار پیچیده‌ای شده است که استفاده از ابزارهای ماشین فرآگیرنده^۱ را در

روش‌های غیرخطی آسان می‌کند [۱۰]. ماشین فرآگیرنده، شاخه‌ای از هوش مصنوعی^۲ است که به طراحی و

کاربرد الگوریتم‌هایی که کامپیوتر را قادر به یادگیری از تجربه می‌سازد، می‌پردازد. از ابزارهای ماشین

فرآگیرنده می‌توان به شبکه عصبی مصنوعی^۳ (ANN) و ماشین بردار پشتیبان^۴ (SVM) اشاره کرد.

^۱ Machine learning

^۲ Artificial intelligence

^۳ Artificial Neural Network

^۴ Support Vector Machine

در بخش اول این پژوهه، کارآیی دو روش مدل‌سازی **MLR** و شبکه عصبی پرسپترون چند لایه^۱ (MLP NN) جهت پیش‌بینی پتانسیل اکسایش-کاهش برخی آنتی‌اکسیدان‌های فنولی با یکدیگر مقایسه شده است. در بخش دوم نیز، عملکرد روش‌های MLP NN، MLR و رگرسیون بردار پشتیبان^۲ (SVR) در مدل‌سازی و پیش‌بینی نیمه‌عمر گروه وسیعی از آلاینده‌های با فنیل‌های چندکلره مورد ارزیابی قرار گرفته است.

^۱ Multi-Layer Perceptron Neural Network

^۲ Support Vector Regression

فصل دوم

سُوری

امروزه از روش‌های پارامتری به طور گسترده‌ای جهت پیش‌بینی رفتار مولکول‌های جدید براساس رفتار مولکول‌های مشابه استفاده می‌شود. در این روش‌ها، بین یک سری توصیف‌کننده‌های مولکولی و فعالیت یا ویژگی مولکول‌های مورد نظر ارتباطی منطقی برقرار می‌شود. توصیف‌کننده‌های مولکولی، پارامترهایی هستند که جنبه‌های مختلف ساختاری مولکول را به صورت کمی نشان می‌دهند. پس از بیان خصوصیات ساختاری مولکول به صورت عدد، بین ساختار و فعالیت مولکول مورد بررسی رابطه‌ای ریاضی یا کمی برقرار می‌شود. به عبارت دیگر در این گونه مطالعات، توصیف‌کننده‌ها به عنوان متغیرهای مستقل و فعالیت یا ویژگی (پارامتر بیولوژیکی یا شیمیایی) مورد نظر به عنوان متغیر وابسته در نظر گرفته می‌شود. این رابطه می‌تواند برای پیش‌بینی پاسخ بیولوژیکی یا شیمیایی دیگر ساختارهای مشابه مورد استفاده قرار گیرد.

در مرحله مدل‌سازی، مدلی از متغیر وابسته بر حسب متغیرهای مستقل ساخته شده و در مرحله پیش‌بینی مورد ارزیابی قرار می‌گیرد.

مراحل کلی مدل‌سازی به روش پارامتری به شرح زیر است:

- جمع‌آوری و انتخاب سری داده‌ها
- رسم و بهینه‌سازی ساختارهای مولکولی
- محاسبه توصیف‌کننده‌های مولکولی
- تجزیه و تحلیل آماری توصیف‌کننده‌ها و انتخاب مؤثرترین آن‌ها
- ایجاد مدل‌های آماری
- تحلیل و ارزیابی آماری مدل‌ها و انتخاب بهترین مدل

۱-۲- جمع‌آوری و انتخاب سری داده‌ها

اولین مرحله مدل‌سازی، جمع‌آوری و انتخاب یک سری مولکولی است که مقادیر تجربی ویژگی یا فعالیت مورد نظر آن‌ها صحیح و قابل اعتماد باشد. انتخاب سری داده‌ها از طریق جستجو در مجلات علمی و پایگاه داده‌ها صورت می‌گیرد. مدل‌سازی باید بر روی ترکیباتی که ساختار مشابهی دارند و فعالیت یا ویژگی آن‌ها در شرایط عملی یکسانی اندازه‌گیری شده است، انجام شود. این امر، لازمه دستیابی به مدلی قابل قبول است تا بتوان نتایج حاصل را به سایر ترکیبات مشابه تعمیم داد. به طور معمول، در مدل‌های خطی و برخی مدل‌های غیرخطی مانند رگرسیون بردار پشتیبان سری داده‌ها به دو گروه سری آموزش^۱ و سری آزمون^۲ تقسیم می‌شود. عملیات مدل‌سازی بر روی سری آموزش که در برگیرنده اکثر مولکول‌ها (۸۰٪) است، انجام می‌گردد. از این‌رو، سری آموزش باید نماینده مناسب و جامعی از خصوصیات ساختاری سری

^۱ Training set

^۲ Test set