

بسم الله الرحمن الرحيم



دانشگاه تربیت مدرس

دانشگاه تربیت مدرس

دانشکده علوم زیستی

پایان نامه کارشناسی ارشد

رشته بیوانفورماتیک

پیش بینی پپتیدهای ضدسرطانی با استفاده از یادگیری ماشین

نگارنده

زینب محمدتبار

استاد راهنما

دکتر پرویز عبدالمالکی

بهمن ۱۴۰۰

چکیده

اگرچه تلاش‌های بسیار زیادی در جهت توسعه درمان‌های جدید سرطان صورت گرفته‌است اما هنوز هم سرطان یکی از علل عمده مرگ در جهان است. شیمی درمانی با وجود عوارض جانبی شدید بر روی سلول‌های طبیعی، هنوز هم روش اصلی مورد استفاده در درمان سرطان در مراحل پیشرفته یا متاستاتیک است. بنابراین، توسعه داروهای ضدسرطانی جدید با سمیت کم برای سلول‌های طبیعی، یک مسیر جدید برای درمان سرطان فراهم می‌کند. در حال حاضر یکی از روش‌های موثر نسبت به شیمی‌درمانی استفاده از پپتیدهای ضدسرطانی است. اما شناسایی این نوع پپتیدها به روش‌های تجربی، پرهزینه و زمانبر است. بنابراین در سال‌های اخیر استفاده از روش‌های محاسباتی جهت شناسایی پپتیدهای ضدسرطانی، مورد توجه پژوهشگران بسیاری قرار گرفته‌است. به‌همین منظور در این پژوهش نیز سعی شده‌است تا توسط یادگیری ماشین مدلی برای پیش‌بینی پپتیدهای ضدسرطانی از پپتیدهای غیرضدسرطانی ایجاد شود.

در این پژوهش با جمع‌آوری داده‌های مربوط به توالی پپتیدهای ضدسرطانی و غیرضدسرطانی، استخراج ویژگی از توالی این پپتیدها توسط کتابخانه پایتون iFeature، ساخت و آموزش دو طبقه‌بند جنگل تصادفی و ماشین بردار پشتیبان توسط کتابخانه پایتون scikit-learn بر روی ویژگی‌های استخراج شده، به نتایج با صحت بالای ۸۲ درصد توسط طبقه‌بند جنگل تصادفی بر روی داده‌های تست مستقل دست یافته‌ایم. همچنین با بررسی نتایج مدل‌های آموزش دیده بر روی داده‌های تست مستقل، می‌توان اعلام کرد که ترکیب دو ویژگی QSOOrder و APseduoAAC توسط هر دو طبقه‌بند جنگل تصادفی و ماشین بردار پشتیبان بالاترین عملکرد برای تشخیص پپتیدهای ضدسرطانی از بین پپتیدهای غیرضدسرطانی را با ۸۱ درصد صحت دارند.

کلید واژه: پیش‌بینی پتیدهای ضدسرطانی، یادگیری ماشین، جنگل تصادفی، ماشین بردار پشتیبان، سرطان

فهرست مطالب

۱ فصل اول: مقدمه	۱
۱-۱ مقدمه‌ای بر پتیدهای ضدسرطانی.....	۱
۲-۱ ویژگی‌های زیستی پتیدهای ضدسرطانی.....	۲
۳-۱ طبقه‌بندی پتیدهای ضدسرطانی	۴
۱-۳-۱ پتیدهای ضدسرطانی با ساختار آلفا-هلیکس (α -Helical ACPs).....	۵
۲-۳-۱ پتیدهای ضدسرطانی غنی از سیستئین.....	۵
۳-۳-۱ پتیدهای ضدسرطانی با ساختار بتا-شیت (β -Sheet ACPs).....	۶
۴-۳-۱ پتیدهای ضدسرطانی غنی از اسیدآمینه‌های منظم.....	۶
۵-۳-۱ پتیدهای ضدسرطانی با اسیدآمینه‌های تغییر یافته یا اصلاح شده.....	۷
۴-۱ روش‌های تجربی شناسایی پتیدهای ضدسرطانی.....	۷
۵-۱ روش‌های محاسباتی پیش‌بینی پتیدهای ضدسرطانی	۸
۱-۵-۱ روش‌های مبتنی بر یادگیری ماشین	۸
۶-۱ پیشینه تحقیق	۱۳
۷-۱ هدف از انجام طرح.....	۱۵
۲ فصل دوم: مواد و روش‌ها.....	۱۷
۲-۱ مجموعه داده.....	۱۷

۱۸مجموعه داده مثبت ۱-۱-۲
۱۹مجموعه داده منفی ۲-۱-۲
۱۹پیش پردازش مجموعه داده ۲-۲
۱۹۲-۲-۱ شناسایی و حذف توالی پپتیدهای مشابه توسط ابزار cd-hit
۲۳۲-۲-۲ فیلتر طول توالی پپتیدها
۲۵۲-۲-۳ بالانس کردن داده‌های مثبت و منفی
۲۷۲-۳ مجموعه داده آموزش و تست
۲۷۲-۳-۱ مجموعه داده آموزش
۲۷۲-۳-۲ مجموعه داده تست مستقل
۲۸۴-۲ استخراج ویژگی
۲۸۲-۴-۱ ویژگی Pseudo-Amino Acid Composition (PseudoAAC)
۳۰۲-۴-۲ ویژگی Amphiphilic Pseudo-Amino Acid Composition (APseudoAAC)
۳۱۲-۴-۳ ویژگی Composition of k-spaced Amino Acid Pairs (CKSAAP)
۳۱۲-۴-۴ ویژگی Composition/Transition/Distribution (CTD)
۳۳۵-۴-۲ ویژگی Dipeptide Deviation from Expected Mean (DDE)
۳۵۶-۴-۲ ویژگی Moran correlation (Moran)
۳۷۲-۴-۷ ویژگی Geary correlation (Geary)

۳۷Normalized Moreau-Broto Autocorrelation (NMBroto) ویژگی ۲-۴-۸
۳۸ k -Spaced Conjoint Triad (KSCTriad) ویژگی ۹-۴-۲
۳۸Quasi-sequence-order (QSOrder) ویژگی ۲-۴-۱۰
۳۹پیش‌پردازش داده ۲-۵
۳۹نرمالسازی داده ۲-۶
۴۰معیارهای ارزیابی ۷-۲
۴۱حساسیت ۲-۷-۱
۴۱اختصاصیت ۲-۷-۲
۴۲صحت ۳-۷-۲
۴۲دقت ۴-۷-۲
۴۲اعتبارسنجی متقابل k -لایه ۲-۸
۴۳انتخاب ویژگی ۲-۹
۴۳روش بسته‌بندی ۲-۹-۱
۴۶کاهش ابعاد مجموعه داده ۱۰-۲
۴۷مدل‌های یادگیری ماشین اجرا شده روی مجموعه داده‌ها ۲-۱۱
۴۸طبقه‌بندی جنگل تصادفی ۱-۱۱-۲
۴۹طبقه‌بندی ماشین بردار پشتیبان ۲-۱۱-۲

۵۱	فصل سوم: نتایج و بحث
۵۱	۳-۱ نتایج طبقه‌بندها
۵۲	۳-۲ نتایج بدست آمده بر روی داده‌های ۱۰ مؤلفه اول PCA
۵۳	۳-۲-۱ نتایج مربوط به طبقه‌بند جنگل تصادفی
۵۶	۳-۲-۲ نتایج مربوط به طبقه‌بند ماشین بردار پشتیبان
۵۹	۳-۳ نتایج ویژگی‌های انتخاب شده توسط الگوریتم انتخاب متوالی رو به جلو
۶۱	۳-۴ نتایج بدست آمده بر روی ۳ مجموعه داده SFS50، SFS100 و SFS200
۶۱	۳-۴-۱ نتایج مربوط به طبقه‌بند جنگل تصادفی
۶۳	فصل چهارم: بحث و پیشنهادات

فهرست شکل ها و نمودارها

- شکل ۱-۱: تعدادی پپتید ضد میکروبی طبیعی با فعالیت ضدسرطانی [۶] ۴
- شکل ۲-۱: ساختار کلی یک درخت تصمیم ۱۰
- شکل ۳-۱: ساختار کلی یک طبقه‌بند جنگل تصادفی ۱۱
- شکل ۴-۱: ساختار کلی یک ماشین بردار پشتیبان دو کلاس خطی ۱۲
- شکل ۵-۱: مثالی از فضای نمونه پس از اعمال تابع هسته ۱۳
- شکل ۶-۱: مراحل ادغام ساخت مجموعه داده از پپتیدهای ضدسرطانی و غیرضدسرطانی [۳۲] ۱۴
- شکل ۱-۲: نمایی کلی از مجموعه داده مثبت (پپتیدهای ضدسرطانی) ۱۸
- شکل ۲-۲: نمایی کلی از مجموعه داده منفی (پپتیدهای غیرضدسرطانی) ۱۹
- شکل ۳-۲: نمایی کلی از مجموعه داده مثبت پس از اعمال cd-hit90 ۲۱
- شکل ۴-۲: نمایی کلی از مجموعه داده منفی پس از اعمال cd-hit90 ۲۱
- شکل ۵-۲: نمایی کلی از مجموعه داده مثبت پس از اعمال cd-hit80 ۲۲
- شکل ۶-۲: نمایی کلی از مجموعه داده منفی پس از اعمال cd-hit80 ۲۲
- شکل ۷-۲: نمایی کلی از مجموعه داده مثبت cd-hit90 پس از اعمال فیلتر طول توالی ۲۴
- شکل ۸-۲: نمایی کلی از مجموعه داده منفی cd-hit90 پس از اعمال فیلتر طول توالی ۲۴
- شکل ۹-۲: نمایی کلی از مجموعه داده منفی cd-hit90 پس از انتخاب رندوم ۲۶
- شکل ۱۰-۲: نمایی کلی از مجموعه داده مثبت و منفی پس از اعمال پیش پردازش ها ۲۶
- شکل ۱۱-۲: یک مثال مصور از ویژگی های فیزیکوشیمیایی در دیتابیس AAindex [۴۴] ۳۵
- شکل ۱۲-۲: نمودار انتخاب متوالی رو به جلو تا ۵۰ ویژگی ۴۵

شکل ۲-۱۳: نمودار انتخاب متوالی رو به جلو تا ۱۰۰ ویژگی..... ۴۵

شکل ۲-۱۴: نمودار انتخاب متوالی رو به جلو تا ۲۰۰ ویژگی..... ۴۶

شکل ۳-۱: نمودار ۵۰ ویژگی انتخاب شده توسط انتخاب متوالی رو به جلو..... ۶۰

شکل ۳-۲: نمودار ۱۰۰ ویژگی انتخاب شده توسط انتخاب متوالی رو به جلو..... ۶۰

شکل ۳-۳: نمودار ۲۰۰ ویژگی انتخاب شده توسط انتخاب متوالی رو به جلو..... ۶۱

فهرست جدول‌ها

- جدول ۳-۱: نتایج بهترین مدل جنگل تصادفی در هر تک‌ویژگی ۵۴
- جدول ۳-۲: نتایج بهترین مدل جنگل تصادفی در هر ترکیب دوتایی ویژگی‌ها ۵۵
- جدول ۳-۳: نتایج بهترین مدل جنگل تصادفی در ترکیب تمام ۷ ویژگی‌ها ۵۶
- جدول ۳-۴: نتایج بهترین مدل ماشین بردار پشتیبان در هر تک‌ویژگی ۵۶
- جدول ۳-۵: نتایج بهترین مدل جنگل تصادفی در هر ترکیب دوتایی ویژگی‌ها ۵۷
- جدول ۳-۶: نتایج بهترین مدل ماشین بردار پشتیبان در ترکیب تمام ۸ ویژگی‌ها ۵۹
- جدول ۳-۷: نتایج بهترین مدل جنگل تصادفی بر روی داده‌های بدست آمده از الگوریتم SFS ۶۲

۱ فصل اول: مقدمه

۱-۱ مقدمه‌ای بر پپتیدهای ضدسرطانی^۱

سرطان یکی از علل عمده مرگ و میر میلیون‌ها انسان در جهان است [۱, ۲]. برآورد آژانس بین‌المللی تحقیقات سرطان (IARC)^۲ نشان می‌دهد که در سال ۲۰۱۸، ۱۸.۱ میلیون مورد جدید سرطان و ۹.۶ میلیون مرگ ناشی از سرطان در سراسر جهان رخ داده‌است. همچنین، آمار در سراسر جهان نشان می‌دهد که سرطان ریه، پستان و روده بزرگ بیش از همه انواع سرطان‌ها بوده‌است [۳].

^۱ Anticancer Peptides (ACP)

^۲ International Agency for Research on Cancer

یکی از روش‌های معمول در درمان سرطان، پرتودرمانی است. اما در روش پرتودرمانی علاوه بر از بین رفتن سلول‌های سرطانی، سلول‌های طبیعی نیز دچار آسیب می‌شوند. مشاهده شده‌است، که سلول‌های سرطانی در برابر داروهای شیمی‌درمانی که جهت مهار تکثیر DNA به کار می‌روند به صورت مقاوم و خاموش به تکثیر خود ادامه می‌دهند[۴]. بنابراین، توسعه داروهای ضدسرطانی جدید با سمیت کم برای سلول‌های طبیعی، یک مسیر جدید برای درمان سرطان فراهم می‌کند. پپتیدها به دلیل عواملی نظیر اندازه کوچک، سنتز راحت، فعالیت و ویژگی بالا و تنوع بیولوژیکی توجه دانشمندان را به خود جلب کرده‌است. از جمله پپتیدهای درمانی که در سال‌های اخیر به منظور درمان سرطان مورد توجه قرار گرفته‌اند، می‌توان پپتیدهای کاتیونی ضدسرطانی و پپتیدهای نفوذپذیر سلولی را نام برد.

تحقیقات نشان می‌دهد که پپتیدهای ضدسرطانی، در مقایسه با روش‌های شیمی‌درمانی و پرتودرمانی، عملکردی مناسب برای شناسایی و مهار سلول‌های سرطانی بدون ایجاد مقاومت دارویی دارند[۵]. بنابراین در حال حاضر یکی از روش‌های موثر در درمان سرطان که توجه دانشمندان را به خود جلب کرده- است، استفاده از پپتیدهای ضدسرطانی است. در ادامه به بررسی ویژگی‌های زیستی پپتیدهای ضدسرطانی می‌پردازیم.

۲-۱ ویژگی‌های زیستی پپتیدهای ضدسرطانی

همانطور که گفته شد یکی از راهکارهای نوین برای درمان سرطان، استفاده از پپتیدهای ضدسرطانی است که با توجه به انتخابی بودن نسبت به سلول‌های سرطانی، عوارض کمتر و اثربخشی بالا مورد توجه بسیاری

از دانشمندان قرار گرفته است [۶]. یکی از انواع پپتیدهای مورد استفاده، پپتید ضد میکروبی^۳ است که بخشی از پاسخ ایمنی ذاتی در برابر میکروب‌ها در بسیاری از گونه‌ها است. این پپتیدها دارای وزن مولکولی کم (۴۰-۱۰ اسیدآمینه) و ساختار آمفی‌پاتیک و کاتیونی هستند که به آن‌ها اجازه می‌دهد تا غشاهای منفی و سلول‌های سرطانی (همانند باکتری‌ها) را هدف قرار دهند. غشاء بیرونی سلول‌های سرطانی نسبت به سلول‌های طبیعی دارای فسفولیپیدهایی با بار منفی بیشتر نظیر فسفاتیدیل سرین و گلیکوپروتئین‌ها و گلیکوز آمینوگلیکان‌های منفی هستند [۷]. پپتیدهای ضدسرطانی با توجه به ویژگی‌های کاتیونی و آمفی‌پاتیک خود با تعاملات الکترواستاتیک به سلول‌های سرطانی متصل شده و از این رو این پپتیدها از طریق نکروز یا آپوپتوز سمیت سلول‌های سرطانی را محدود می‌کنند [۸-۱۰].

پپتیدهای ضدسرطانی در درمان دیابت، سرطان و بیماری‌های قلبی عروقی به کار برده می‌شوند [۱۱]. اکثر پپتیدهای ضدسرطانی به صورت محلول و بدون ساختار هستند و تنها پپتید ضد میکروبی انسانی به نام 37LL با توجه به حضور یک یا چند پیوند دی‌سولفیدی با صفحات β -sheet به شکل گلوبمرولی دیده شده است. از نظر تنوع ساختاری چند صد پپتید ضد میکروبی مورد مطالعه قرار گرفته است. این پپتیدها دارای ساختار آلفاهلیکسی غنی از سیستئین (نظیر Cecropins) یا دارای ساختار β -sheet (نظیر Defensins) هستند، همچنین حضور اسیدآمینه آرژنین، پرولین، هیستیدین و تریپتوفان در این پپتیدها (نظیر Indolicidin) معمول است [۱۲].

³ Antimicrobial Peptides (AMP)

Table 1
Summary of select naturally occurring cationic antimicrobial peptides with anticancer activities

Peptide	Source	Primary amino acid sequence ^a	Class	Net ^b	Anticancer activity
BMAP-28	<i>Bos taurus</i>	GGLRSLGRKILRAWKKYGPIIVPIIRI	α -Helix	+7	Membranolytic
HNP-1 (β -defensin)	<i>Homo sapiens</i>	AC ₁ YC ₂ RPAC ₃ LAGERRYGTC ₂ IYQGRLWAFC ₃ C ₁	β -Sheet	+3	Membranolytic Antiangiogenic?
Lactoferricin B	<i>Bos taurus</i>	FKC ₁ RRWQWRMKKLGAPSTC ₁ VRRAF	β -Sheet	+8	Membranolytic Apoptosis inducer
LL-37	<i>Homo sapiens</i>	LLGDFFRKSKKEKIGKEFKRIVQRIKDFLRNLPRTES	α -Helix	+6	Antiangiogenic
Magainin 2	<i>Xenopus laevis</i>	GIGKFLHSAKKFGKAFVGEIMNS	α -Helix	+3	Membranolytic
Melittin	<i>Apis mellifera</i>	GIGAVLKVLTTGLPALISWIKRKRQQ	α -Helix	+6	Apoptosis inducer? Membranolytic
Tachyplesin I	<i>Tachyplesus tridentatus</i>	KWC ₁ FRVC ₂ YRGIC ₂ YRRC ₁ R	β -Sheet	+6	PLA ₂ ^c activator PLD ^d activator Binds hyaluronan and activates complement (C1q) Antiangiogenic? Induces cancer cell differentiation

^a Amino acid sequences are given in one-letter code. Subscripts indicate pairings of Cys residues that form disulfide bonds. Boldface indicates cationic amino acid residues.

^b At neutral pH.

^c Phospholipase A₂.

^d Phospholipase D.

شکل ۱-۱: تعدادی پپتید ضد میکروبی طبیعی با فعالیت ضد سرطانی [۶]

در دهه‌های اخیر استفاده از این پپتیدها به عنوان عوامل ضد سرطان به عنوان یک روش درمانی نوین در نظر گرفته شده است [۱۳-۱۵].

۳-۱ طبقه‌بندی پپتیدهای ضد سرطانی

رزونانس مغناطیسی هسته‌ای ^۴NMR به عنوان یک روش مفید برای مطالعه جزئیات ساختاری بسیاری از پپتیدهای ضد میکروبی شناخته شده، پدید آمده است. تجزیه و تحلیل ساختار سه بعدی این پپتیدها منجر به درک بهتری از عملکرد این پپتیدها می‌گردد [۱۶]. براساس ساختارهای شناخته شده توسط NMR پپتیدهای ضد سرطانی به پنج گروه طبقه‌بندی می‌شوند.

^۴ NMR spectroscopy

۱-۳-۱ پپتیدهای ضدسرطانی با ساختار آلفا-هلیکس (α -Helical ACPs)

برای اولین بار در پپتید Cecropins بسیاری از ویژگی‌ها شناسایی شد و مطالعات اولیه با NMR نشان داد که پپتید Cecropin-A دارای ساختار هلیکس است. در ساختار هلیکس این پپتید ۱۵٪ الکل Hexafluoro Isopropyl شرکت دارد. در نتیجه پپتیدهای ضد میکروبی دارای ساختار هلیکس، پپتیدهایی آمفی‌پاتیک با سطوح آبگریز و حاوی بار خالص مثبت هستند. پپتیدهای ضد میکروبی با ساختار هلیکس بیشتر از سایر ساختارها مشاهده می‌شوند. پپتید ضد میکروبی به نام Magainins گروه دیگری از پپتیدها با ساختار آلفا-هلیکس است. این پپتیدها از پوست نوعی قورباغه آفریقایی جدا می‌شوند. مطالعات NMR نشان داده‌است، که Cecropins نیز مانند Magainins دارای ساختار آمفی‌پاتیک-هلیکس در ۲۵٪ از Trifluoroethanol است [۱۷].

۱-۳-۲ پپتیدهای ضدسرطانی غنی از سیستئین

پپتیدهای نوتروفیل انسان HNP-1-2-3 اولین پپتیدهای غنی از سیستئین جدا شده از گرانول انسان بودند [۱۸]. پپتیدهای غنی از سیستئین با ۳۰ اسید آمینه در طیف گسترده‌ای از موجودات زنده وجود دارد. این پپتیدها به صورت یک موتیف حفاظت شده حاوی شش سیستئین با سه پیوند دی‌سولفید درون مولکولی است. اکثراً موقعیت پل دی‌سولفیدی بین C1-C4، C2-C3-C5-C6 است. مطالعات کریستالوگرافی HNP-3 توسط XRAY، در ترکیب با سانتریفوژ تعادل رسوب، نشان می‌دهد که این پپتید به صورت دایمر وجود دارد [۱۹]. در مطالعات NMR بر روی ساختار Defensin مشخص گردید، که این پپتید حاوی سه رشته آنتی‌پارالل است. پپتید Drosomycin جدا شده از مگس سرکه شامل چهار پیوندهای دی‌سولفید است و از سه رشته آنتی‌پارالل ساخته شده‌است. در این پپتید یک ساختار هلیکس بین دو رشته اول وجود دارد [۲۰].

۳-۳-۱ پپتیدهای ضدسرطانی با ساختار بتا-شیت (β -Sheet ACPs)

تعدادی از پپتیدهای ضدسرطانی با طول ۲۰ اسیدآمینه دارای یک ساختار سنجاق سری هستند. این پپتیدها دارای یک یا دو پیوند دی‌سولفیدی می‌باشند. پپتیدهای خرچنگ نعل اسبی، Tachyplesin و Polyphemusin II که به صورت موتیف کوتاه سنجاق سری هستند، این پپتید کوتاه توسط دو پیوند دی‌سولفیدی پایدار می‌گردد [۲۰، ۲۱]. مطالعات NMR همراه با ساختارهای 3D نشان می‌دهد که پپتید Tachyplesin شباهت بسیار زیادی با پپتیدهای جدا شده از خوک دارد. مولکول پپتیدی صفحات آنتی پارالل به یک ساختار Turn متصل شده‌اند و شامل دو پل دی‌سولفید هستند [۲۱].

۴-۳-۱ پپتیدهای ضدسرطانی غنی از اسیدآمینه‌های منظم

تعدادی از پپتیدهای ضدسرطانی دارای تعداد زیادی از یک اسیدآمینه خاص هستند. این پپتیدها دارای ساختارهای مختلف آلفا-هلیکس و بتا-شیت هستند. پپتید Histatin، غنی از اسیدآمینه هیستیدین از بزاق انسان جدا شده و این پپتید در برابر کاندید آلبیکانس فعال است [۲۲]. در حالی که پپتیدهای Cathelicidins، غنی از پرولین و دارای ساختار نامنظم هستند. پپتیدهای Tritrpticin سرشار از تریپتوفان و پپتیدهای Bactenecin Bac-5 and Bac-7 مانند پپتیدهای Cathelicidins غنی از اسیدآمینه پرولین هستند؛ در حالی که پپتید PR-39، غنی از اسیدآمینه آرژنین است [۲۳-۲۶].

۱-۳-۵ پپتیدهای ضدسرطانی با اسیدآمینه‌های تغییر یافته یا اصلاح شده

تعدادی از پپتیدهای ضدسرطانی دارای اسیدآمینه‌های اصلاح شده هستند. بهترین نمونه از این پپتیدها، پپتیدهای تولید شده توسط باکتری‌ها است. از این پپتیدها می‌توان به نیسین^۵، که توسط باکتری لاکتوکوکوس لاکتیس^۶ تولید می‌شود، اشاره کرد. این پپتید باکتریایی از اسیدآمینه‌های نادر مانند (3-Methyllationine, Dyhydroalanine and Dehydrobutyrine) تشکیل شده است و این پپتیدها بر علیه باکتری‌های گرم مثبت فعال هستند [۲۷]. از دیگر پپتیدهای ضد میکروبی که دارای اسیدآمینه اصلاح شده هستند، می‌توان به پپتید Leucocin A، که از *Leuconostoc gelidum* جدا شده است، اشاره کرد. این پپتید ساختار کانفورماسیونی آمفی پاتیک از خود نشان می‌دهد، که این ساختارها دارای نقش مهمی در تعامل با غشا هست [۲۸].

۱-۴ روش‌های تجربی شناسایی پپتیدهای ضدسرطانی

برای شناسایی پروتئین‌ها یا پپتیدها به روش تجربی اول از همه نیازمند استخراج نمونه آزمایشگاهی و یا آزمایش‌هایی روی موجود زنده است. همچنین برای کشف عملکرد مورد نظر که در این پژوهش خاصیت ضدسرطانی پپتیدها مدنظر است نیازمند روش‌های آزمایشگاهی خاص خود است. تمام مراحل از یافت نمونه تا بررسی عملکرد پپتید نیازمند متخصص و تجهیزات زیستی است. این پر واضح است که انجام روش‌های آزمایشگاهی برای شناسایی عملکرد یک پپتید نیازمند زمان و هزینه مالی و انسانی قابل توجهی است. به همین دلیل همیشه امید است که با مطالعه بر روی روش‌های محاسباتی بتوان محققان زیستی را تا حدی

⁵ Nissen

⁶ Lactococcus Lactis

نجات داد، تا محققان زیستی برای انجام آزمایش‌های خود، از بین کاندیدهای میلیونی برای مطالعه و آزمایش به کاندیدهای صدتایی یا حتی هزارتایی برسند. این پدیده باعث میشود تا محققان زیستی در زمان و هزینه‌های مالی و مواد صرفه‌جویی کنند و این زمان را صرف تحقیقات با ارزش‌تری کنند. شناسایی پتیدهای ضدسرطانی نیز از این قاعده مستثنی نمی‌شود. به‌همین دلیل در این پایان‌نامه نیز سعی می‌شود تا با انجام یک روش محاسباتی به پیش‌بینی‌ای از پتیدهای ضدسرطانی برسیم.

۱-۵ روش‌های محاسباتی پیش‌بینی پتیدهای ضدسرطانی

همانطور که در قسمت قبل اشاره شد، روش‌های تجربی نیازمند صرف انرژی انسانی و هزینه‌های مالی زیاد است، در نتیجه نیاز به رویکردهای محاسباتی در کنار روش‌های تجربی وجود دارد. یکی از روش‌های محاسباتی برای پیش‌بینی پتیدهای ضدسرطانی، استفاده از داده‌های بدست آمده توسط محققان زیستی برای یادگیری ماشین است. در ادامه به توضیح مختصری از الگوریتم‌های یادگیری ماشین به کاررفته در این پایان‌نامه می‌پردازیم.

۱-۵-۱ روش‌های مبتنی بر یادگیری ماشین

یادگیری ماشین به مجموعه وسیعی از روش‌های محاسباتی برای درک داده‌ها اشاره دارد که می‌تواند الگوهای پیچیده در یک مجموعه را پیدا و بر اساس آن‌ها تصمیم‌گیری کند. این روش‌ها را می‌توان به دو دسته یادگیری با نظارت^۷ و یادگیری بدون نظارت^۸ تقسیم کرد. یادگیری ماشین با نظارت، یک مدل آماری

^۷ Supervised

^۸ Unsupervised

برای پیش‌بینی برچسب کلاس^۹ بر اساس یک یا چند توصیف‌کننده^{۱۰} (ویژگی^{۱۱}) ایجاد می‌کند. در این روش با هدف پیش‌بینی دقیق کلاس‌ها برای مشاهدات آینده، مدلی^{۱۲} تنظیم می‌شود که برچسب کلاس‌ها را به ویژگی‌ها مرتبط کند.

در یادگیری ماشین بدون نظارت، به توصیف‌کننده‌ها کلاسی نسبت داده نشده‌است. در این روش با پیدا کردن ساختار و روابط بین توصیف‌کننده‌ها، داده‌ها کلاستر می‌شوند.^{۱۳} در مسائلی که داده‌ها برچسب کلاس دارند؛ از روش‌های با نظارت که به آن‌ها طبقه‌بندی^{۱۴} می‌گویند، استفاده می‌شود. در ادامه به اختصار طبقه‌بندهای استفاده شده در این پایان‌نامه توضیح داده می‌شود.

۱-۱-۵-۱ طبقه‌بند درخت تصمیم^{۱۵}

درخت نوع خاصی از گراف است و درخت تصمیم یک روش برای پیش‌بینی بر اساس این ساختار است. به دلیل تفسیرپذیر بودن درخت تصمیم، این روش از پرکاربردترین مدل‌ها در مسائل یادگیری ماشین کلاسیک است. این ساختار از گره ریشه^{۱۶} شروع می‌شود و به برگ‌ها^{۱۷} (گره پایانی) ختم می‌شود. برگ است که مشخص می‌کند یک نمونه در چه کلاسی قرار خواهد گرفت. گره‌های داخلی با استفاده از قوانینی و با توجه به ویژگی‌ها، نمونه‌ها را به زیر مجموعه‌های مختلف تقسیم می‌کند، تا در نهایت مشخص شود هر نمونه به کدام کلاس تعلق دارد [۲۹].

⁹ Class Label

¹⁰ Descriptor

¹¹ Feature

¹² Predictors

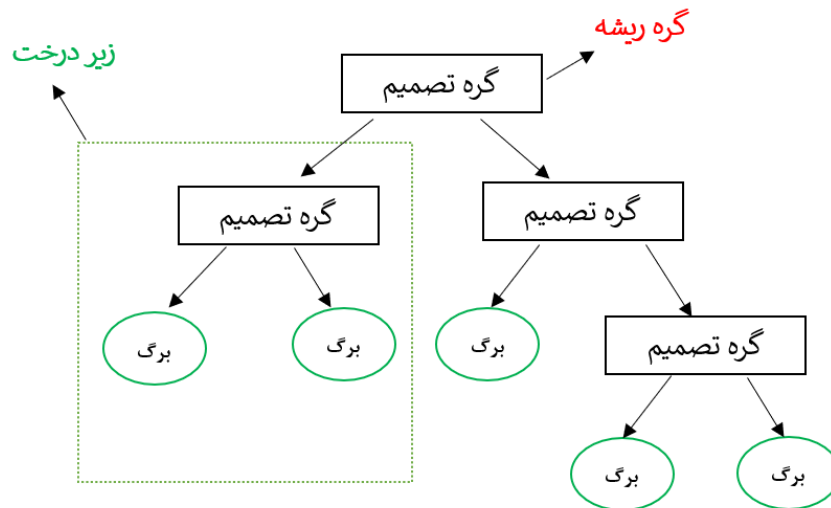
¹³ Clustering

¹⁴ Classification

¹⁵ Decision Tree Classifier

¹⁶ Root Node

¹⁷ Leaves



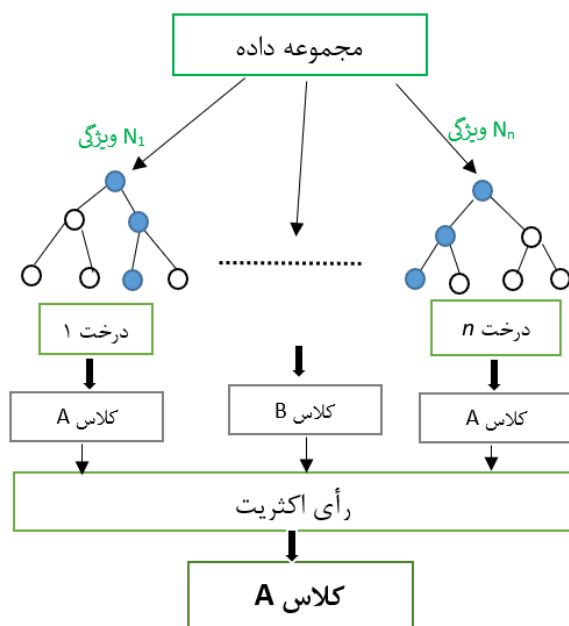
شکل ۱-۲: ساختار کلی یک درخت تصمیم

۲-۱-۵-۱ طبقه‌بند جنگل تصادفی^{۱۸}

طبقه‌بند جنگل تصادفی از تعداد زیادی درخت تصمیم استفاده می‌کند. به همین دلیل الگوریتم جنگل-تصادفی در دسته الگوریتم‌های یادگیری گروهی^{۱۹} قرار می‌گیرد. در این روش هر درخت تصمیم به صورت تصادفی تعدادی از ویژگی‌ها را انتخاب می‌کند و در مورد هر نمونه پیش‌بینی انجام می‌دهد. برای تعیین برچسب یک نمونه بین تمام برچسب‌هایی که هر درخت به آن نمونه نسبت داده است؛ رأی اکثریت گرفته می‌شود [۳۰].

¹⁸ Random Forest Classifier

¹⁹ Ensemble Learning



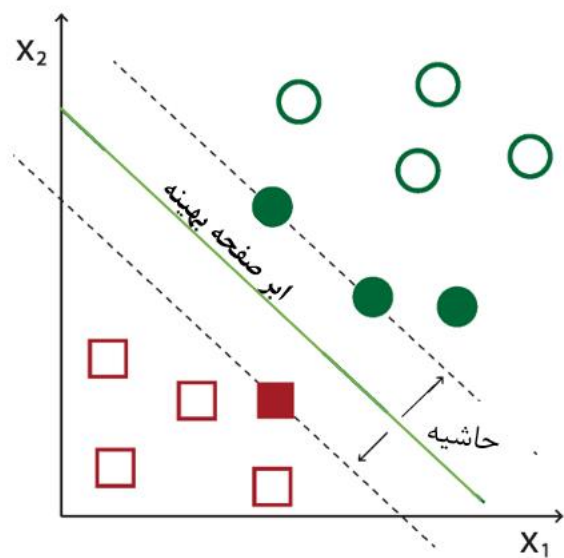
شکل ۳-۱: ساختار کلی یک طبقه‌بند جنگل تصادفی

۳-۱-۵-۱ طبقه‌بند ماشین بردار پشتیبان^{۲۰}

ماشین بردار پشتیبان یک الگوریتم یادگیری ماشین است که سعی دارد با ایجاد یک ابرصفحه داده‌ها را به نسبت تعداد برچسب‌ها تقسیم کند. این طبقه‌بند در ایجاد ابرصفحه مناسب به دنبال بیشینه کردن فاصله نمونه‌ها از ابرصفحه است. به این فاصله حاشیه^{۲۱} گفته می‌شود.

²⁰ Support Vector Machine

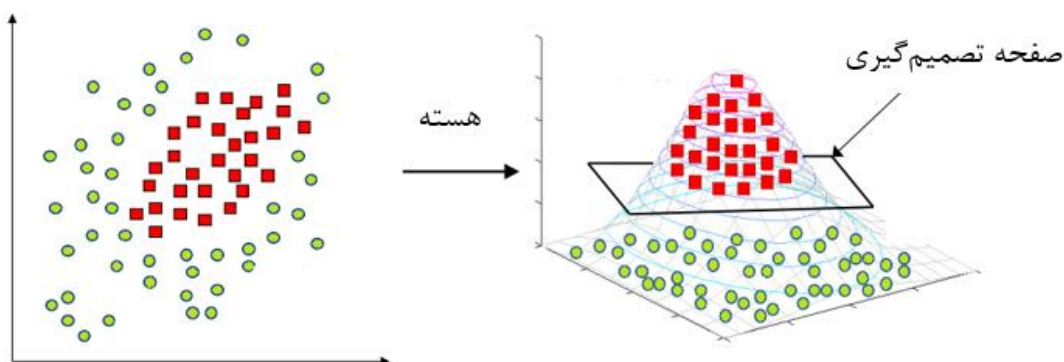
²¹ Margin



شکل ۱-۴: ساختار کلی یک ماشین بردار پشتیبان دوکلاس خطی

از آنجایی که همیشه و به راحتی امکان یافتن چنین ابرصفحه‌ای نیست؛ با تابعی به نام هسته^{۲۲} نمونه‌ها به فضایی جدید منتقل می‌شوند و سپس ابرصفحه مناسب ایجاد می‌شود. همچنین این الگوریتم اجازه می‌دهد تعداد کمی از نمونه‌ها در سمت نادرست ابرصفحه قرار بگیرند [۳۱].

²² Kernel



شکل ۱-۵: مثالی از فضای نمونه پس از اعمال تابع هسته

۶-۱ پیشینه تحقیق

پایگاه داده‌های متعددی برای ذخیره داده پپتیدهای ضد میکروبی با عملکرد ضد سرطانی که توسط محققان زیستی گزارش شده‌اند وجود آمده‌است. همچنین مقالات متعددی در حوزه پیش‌بینی پپتیدهای ضد سرطانی به روش‌های محاسباتی در حال انتشار هستند که این نشان‌دهنده جذابیت و اهمیت موضوع تشخیص پپتیدهای ضد سرطانی از پپتیدهای غیر ضد سرطانی است.

همچنین الگوریتم‌های یادگیری ماشین کلاسیک و مدرن (شبکه‌های عصبی عمیق^{۲۳}) متعددی تا به امروز ارائه شده‌است. یکی از روش‌های پیش‌بینی ACPs^{۲۴} استخراج ویژگی از توالی این پپتیدها است. مطالعات گسترده‌ای در حوزه استخراج ویژگی از توالی پپتیدها به صورت محاسباتی صورت گرفته‌است، که شامل استخراج ویژگی از توالی ساختار اول، ساختار دوم و سوم آن‌ها می‌شود. عمده تحقیقات برای پیش‌بینی

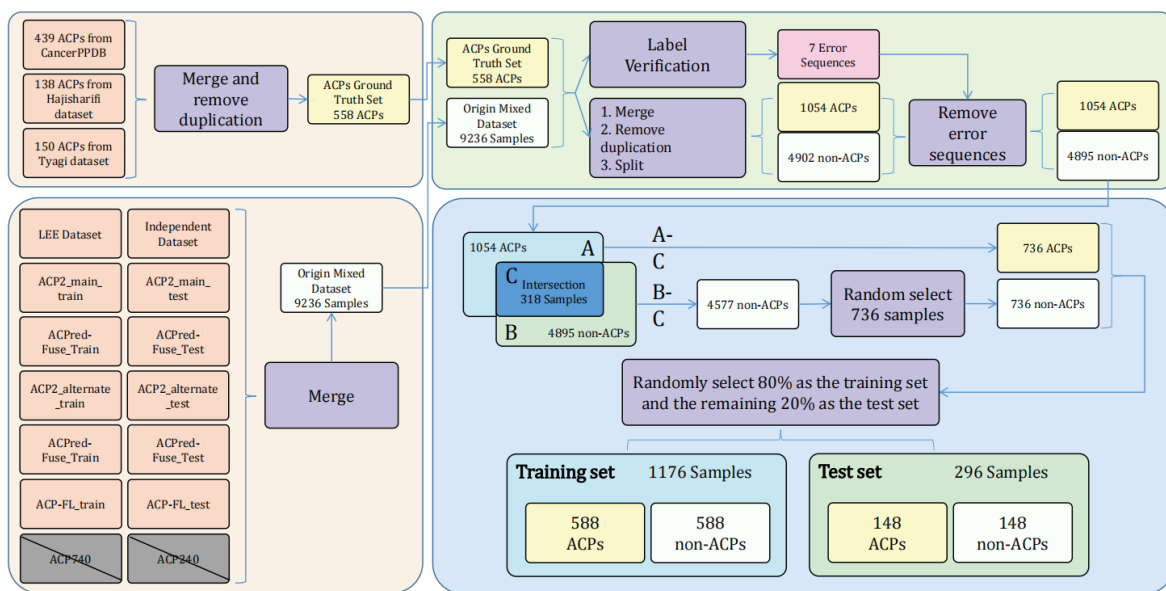
²³ Deep Neural Networks

²⁴ Anticancer Peptides

ACPها توسط استخراج ویژگی از ساختار اولیه این پپتیدها انجام شده است. بعلاوه ابزارهای متعددی برای استخراج این ویژگیها در حال توسعه هستند و منتشر شده اند.

با مطالعه و جستجو بین پایگاه داده ها^{۲۵} و مقالات منتشر شده در حوزه پیش بینی ACPs، مطالعاتی در سال ۲۰۲۱ تجمیعی از تمام داده های موجود در حوزه پپتیدهای ضدسرطانی را در مقاله خود انجام داده است، داده های این مطالعه مورد استفاده این پایان نامه قرار گرفته است [۳۲].

مراحل تجمیع داده بین پایگاه داده های پپتیدهای ضد میکروبی و داده های استفاده شده در مقالات پیشین حوزه ACPs در شکل ۱-۶ مشاهده می شود [۳۲].



شکل ۱-۶: مراحل ادغام ساخت مجموعه داده از پپتیدهای ضدسرطانی و غیر ضدسرطانی [۳۲]

برای استخراج ویژگی از ساختار اولیه، ثانویه و یا ساختار سوم پروتئین و پپتید، ابزارهای متعددی منتشر شده و در حال توسعه هستند. با بررسی ابزارهای متفاوت برای استخراج ویژگی‌ها از ساختار اولیه توالی پروتئین و پپتیدها، پکیج پایتون iFeature که در سال ۲۰۱۸ منتشر شده است [۳۳]، پوشش دهنده نیازهای این پایان‌نامه است.

همچنین با بررسی مقالات موجود در حوزه ACP prediction به روش‌های محاسباتی و همچنین بررسی الگوریتم‌های کلاسیک یادگیری ماشین، بررسی شده‌است که RF^{۲۶} و SVM^{۲۷} کاندیدای خوبی برای استفاده بر روی داده‌های ACP، non-ACP^{۲۸} هستند [۳۴-۳۹].

۷-۱ هدف از انجام طرح

اگرچه AMPs^{۲۹} در چندین دهه قبل شناخته شده‌اند، اما تنها در دهه اخیر است که تعداد مقالات مربوط به فعالیت‌های ضدسرطانی آن‌ها افزایش یافته و از آن‌ها به عنوان پپتیدهای ضدسرطانی یاد میکنند [۶-۸، ۱۴، ۴۰]. به همین علت بر این باوریم که در سال‌های آینده، این پپتیدها به علت ویژگی‌های منحصر به فردشان در راستای تأثیرگذاری روی سلولهای سرطانی، پیشرفت مهمی در درمان بیماری سرطان که از بزرگترین نگرانی‌های جامعه بشری در جهان است، رقم خواهند زد. استراتژی دیگری که مورد توجه قرار گرفته‌است، استفاده ترکیبی از پپتیدها با داروهای مرسوم شیمی درمانی است که هزینه‌های درمان را کاهش می‌دهد و باعث به حداقل رساندن مشکل مقاومت به سرطان و جلوگیری از عود مجدد آن میشود. پیشرفت‌هایی در جهت تولید این پپتیدها در مقیاس بزرگ در جهان صورت گرفته‌است تا این روش درمانی، برای

²⁶ Random Forest

²⁷ Support Vector Machine

²⁸ Non Anticancer Peptide

²⁹ Antimicrobial Peptides

بیماران ارزان تر و قابل دسترس تر باشد. هرچند محدودیت‌هایی ازجمله شباهت احتمالی این پتیدها با آنتی‌ژنهای خودی یا تحریک سیستم ایمنی علیه این پتیدها میتواند وجود داشته باشد. در نهایت با توجه به مطالب گفته‌شده میتوان پیش‌بینی کرد این پتیدها مسیری رو به پیشرفت در جهت بهینه‌سازی روند درمان بیماری سرطان را طی کرده و میتوانند یک روش درمانی نوین و با عوارض کم را ارائه دهند. همچنین با توجه به هزینه مالی و انسانی زیاد برای شناسایی این پتیدها به روش آزمایشگاهی، همیشه نیاز به روش‌هایی محاسباتی در کنار روش‌های آزمایشگاهی وجود دارد. به همین دلیل در این پایان‌نامه نیز سعی در بررسی پیش‌بینی پتیدهای ضدسرطانی از بین پتیدهای غیرضدسرطانی توسط یادگیری ماشین شده- است.

۲ فصل دوم: مواد و روش‌ها

۲-۱ مجموعه داده^{۳۰}

مجموعه داده مورد استفاده برای پیش‌بینی پپتیدهای ضدسرطانی که در این پایان‌نامه استفاده شده و در فصل قبل نیز به آن‌ها اشاره شد به دو مجموعه تقسیم می‌شوند [۳۲].

۱- مجموعه داده مثبت^{۳۱}، شامل پپتیدهای ضد میکروبی با عملکرد ضدسرطانی، که توسط محققان زیستی به صورت آزمایشگاهی تایید شده‌اند.

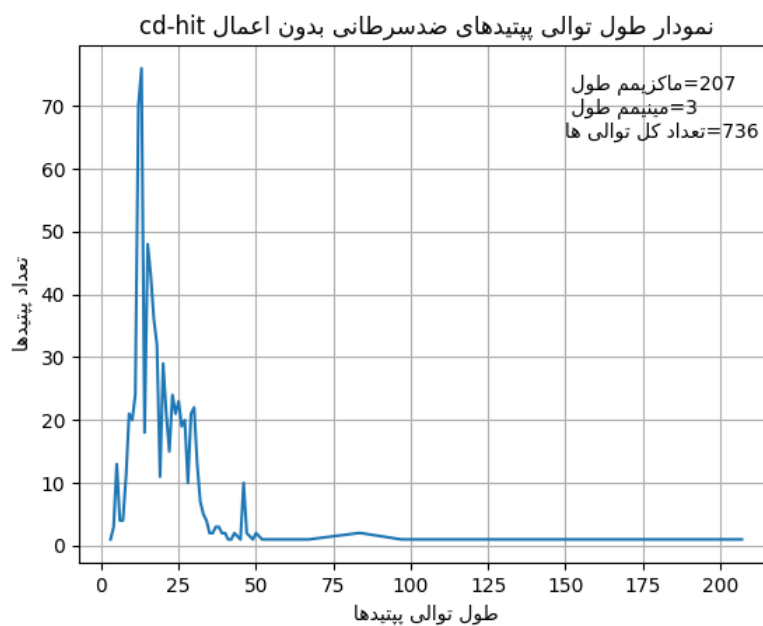
³⁰ Datasets

³¹ Positive Dataset

۲- مجموعه داده منفی^{۳۲}، شامل پپتیدهای ضد میکروبی که عملکرد ضدسرطانی برای آنها گزارش نشده است.

۲-۱-۱ مجموعه داده مثبت

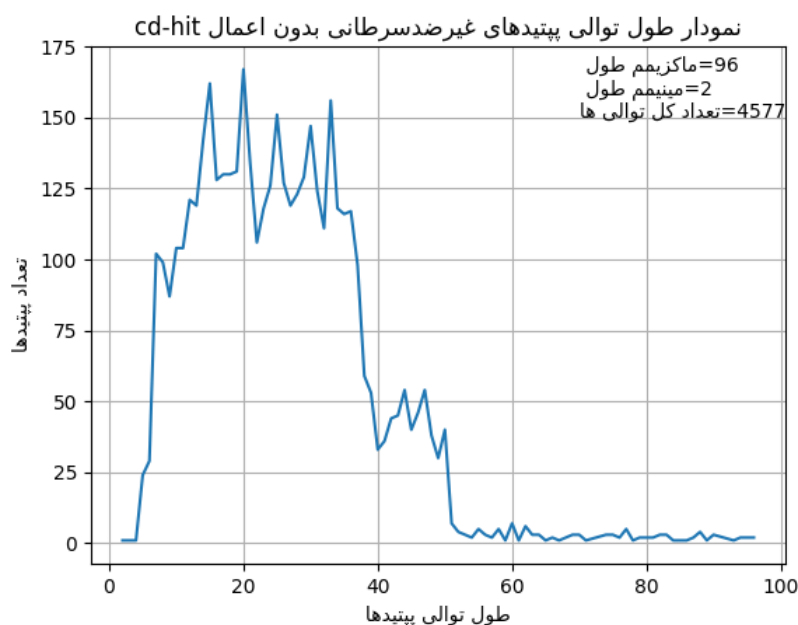
مجموعه داده مثبت شامل پپتیدهای ضد میکروبی با عملکرد ضدسرطانی، که توسط محققان زیستی به صورت آزمایشگاهی تایید شده‌اند، شامل ۷۳۶ پپتید ضدسرطانی (ACP) است [۳۲].



شکل ۲-۱: نمای کلی از مجموعه داده مثبت (پپتیدهای ضدسرطانی)

۲-۱-۲ مجموعه داده منفی

مجموعه داده منفی شامل پپتیدهای ضد میکروبی که توسط محققان زیستی، عملکرد ضد سرطانی برای آنها گزارش نشده است، شامل ۴۵۷۷ پپتید غیر ضد سرطانی (non-ACP) است [۳۲].



شکل ۲-۲: نمایی کلی از مجموعه داده منفی (پپتیدهای غیر ضد سرطانی)

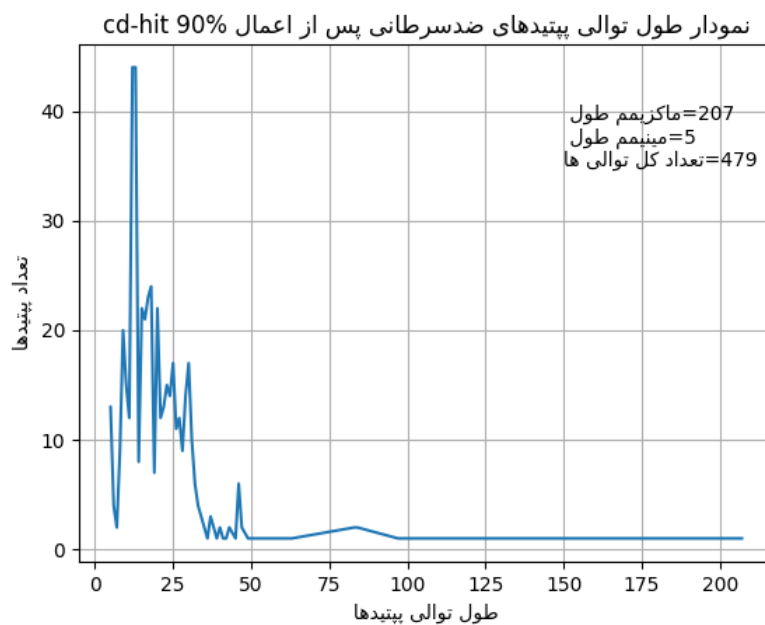
۲-۲ پیش پردازش مجموعه داده

۲-۲-۱ شناسایی و حذف توالی پپتیدهای مشابه توسط ابزار cd-hit

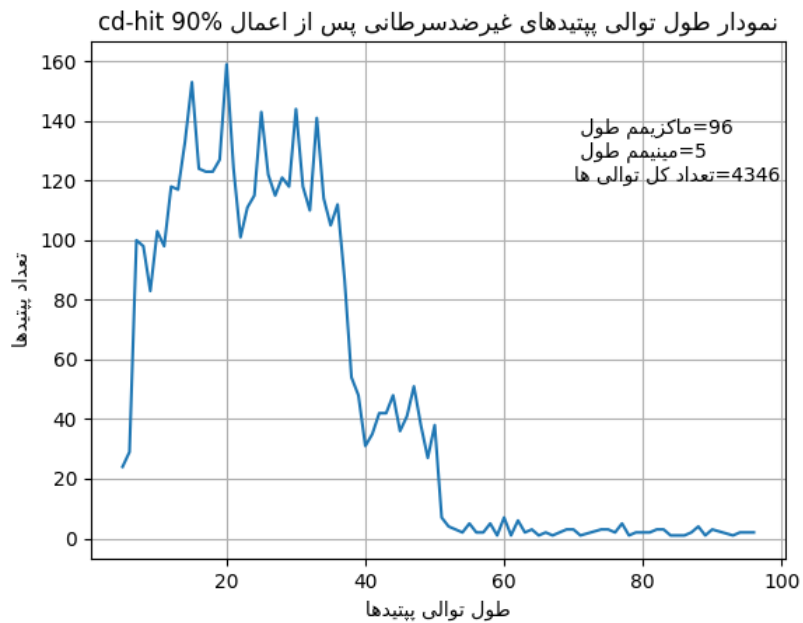
اجتماع مجموعه داده خام مثبت و مجموعه داده خام منفی، مجموعاً تعداد ۵۳۱۳ پپتید است. با استفاده از ابزار cd-hit [۴۱]، ۲ بار بصورت کاملاً مستقل اقدام به حذف توالی‌های مشابه شد. یکبار بر روی داده‌های

خام cd-hit با ۹۰ درصد شباهت^{۳۳} و حذف پپتیدهای با طول کمتر از ۵ اسیدآمینو اجرا شد. بار دیگر مجدداً بر روی داده‌های خام cd-hit با ۸۰ درصد شباهت و حذف پپتیدهای با طول کمتر از ۵ اسیدآمینو اجرا شد. نتایج بدست آمده به تفکیک داده‌های مثبت و منفی در شکل‌های زیر مشاهده می‌شود. در نهایت با اجرا cd-hit90 مجموع داده به ۴۸۲۵ داده کاهش یافت و با اجرا cd-hit80 مجموع داده به ۴۵۳۴ داده کاهش یافت. بدلیل اینکه با اعمال cd-hit80 مقدار قابل توجهی از داده‌ها کم شد، داده‌های مربوط به cd-hit90 ملاک انجام ادامه پژوهش قرار گرفتند.

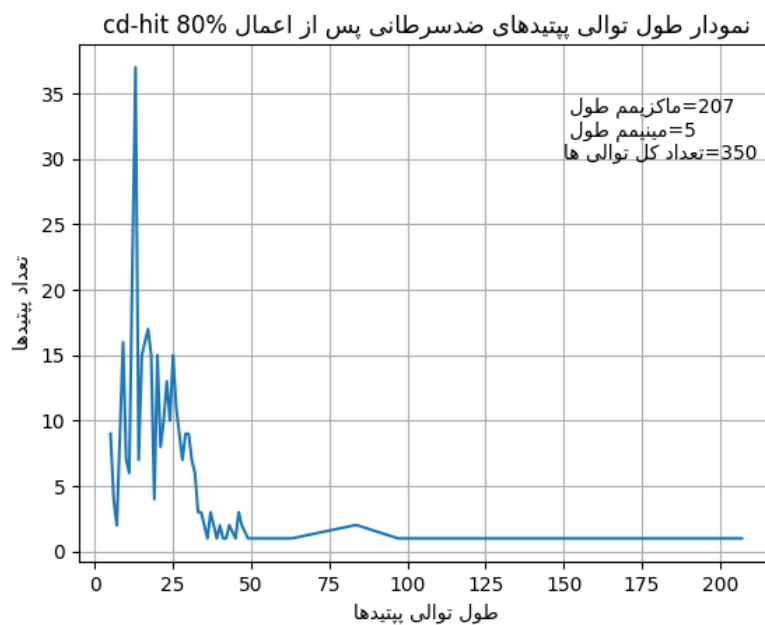
³³ Similarity



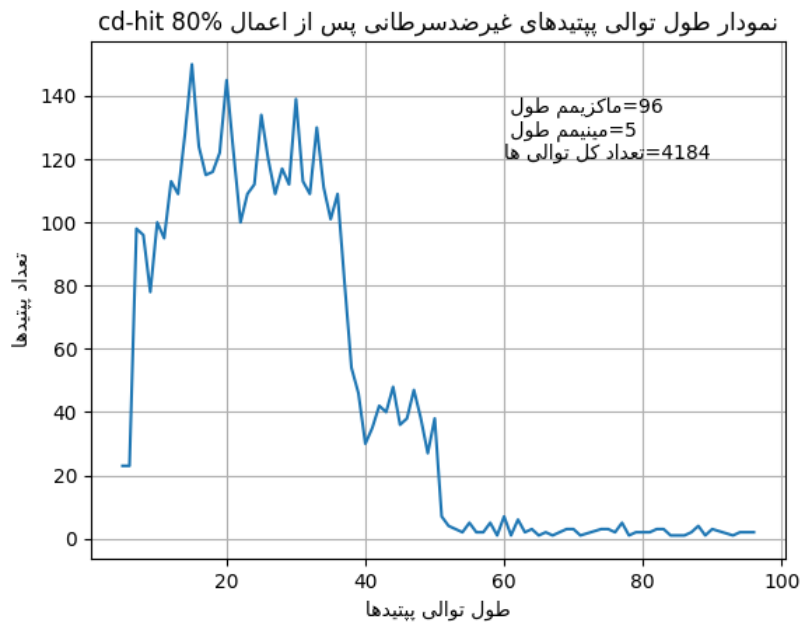
شکل ۲-۳: نمایی کلی از مجموعه داده مثبت پس از اعمال cd-hit90



شکل ۲-۴: نمایی کلی از مجموعه داده منفی پس از اعمال cd-hit90



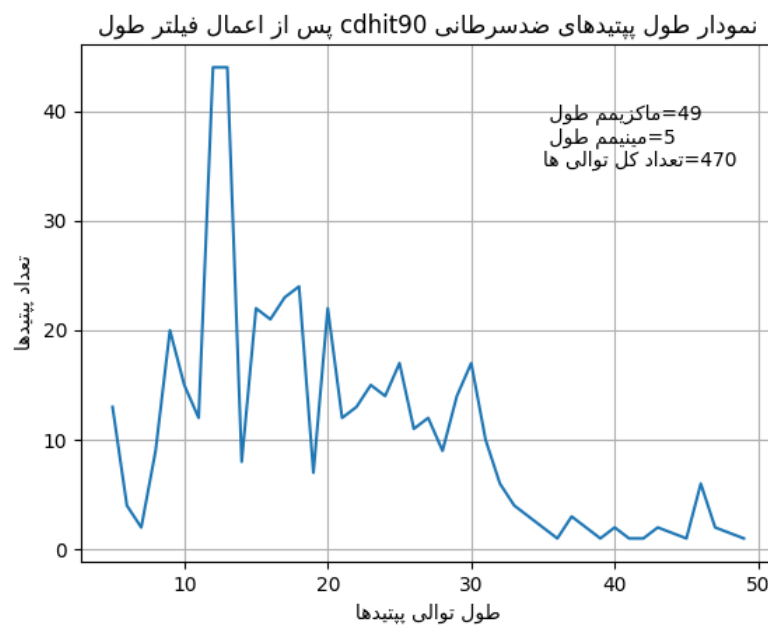
شکل ۲-۵: نمایی کلی از مجموعه داده مثبت پس از اعمال cd-hit80



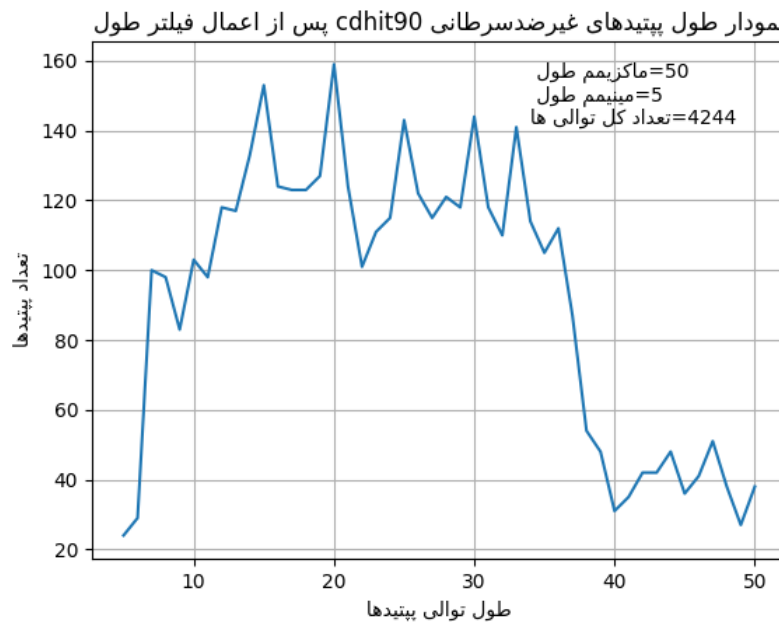
شکل ۲-۶: نمایی کلی از مجموعه داده منفی پس از اعمال cd-hit80

۲-۲-۲ فیلتر طول توالی پپتیدها

با توجه به اینکه نرم طول پپتیدهای ضد میکروبی در طبیعت بین ۵ تا حداکثر ۵۰ اسید آمینه دارند [۶]، به همین دلیل بر روی داده‌های cdhit-90 فیلتر با طول حداکثر ۵۰ اعمال شد، که مجموعه داده به ۴۷۱۴ داده کاهش یافت.



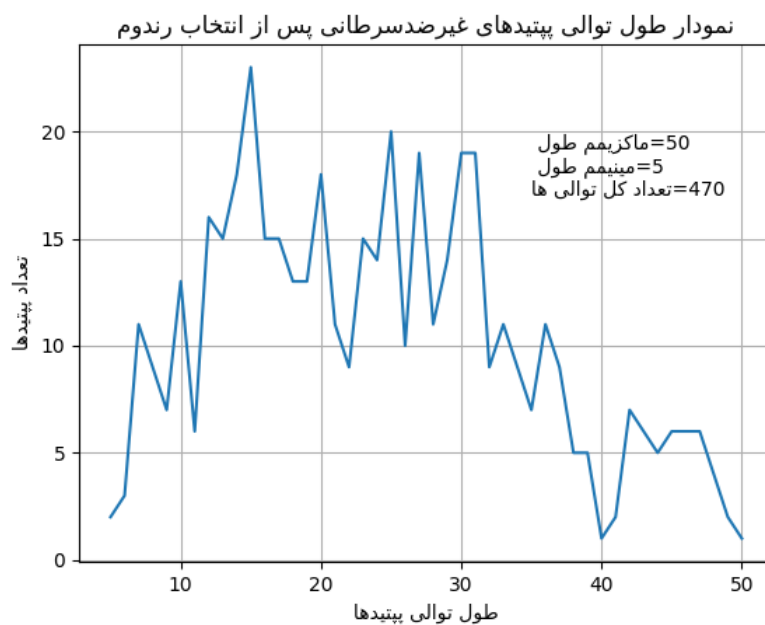
شکل ۲-۷: نمایی کلی از مجموعه داده مثبت cd-hit90 پس از اعمال فیلتر طول توالی



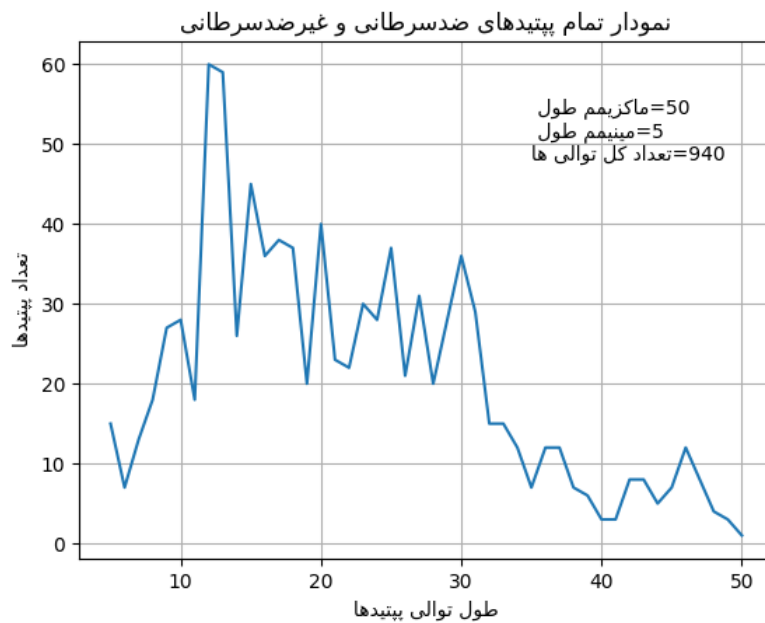
شکل ۲-۸: نمایی کلی از مجموعه داده منفی cd-hit90 پس از اعمال فیلتر طول توالی

۲-۲-۳ بالانس کردن داده‌های مثبت و منفی

پس از اعمال پیش‌پردازش بر روی داده‌های خام، تعداد داده‌های cd-hit90 به تعداد ۴۷۰ داده مثبت (ACP) و ۴۲۴۴ داده منفی (non-ACP) کاهش یافت. با توجه به تعداد بیشتر مجموعه داده منفی و همچنین عدم گزارش وسیع پپتیدهای ضدسرطانی (ACP) به صورت آزمایشگاهی، یکی از چالش‌ها انتخاب مجموعه داده منفی است. برای انتخاب داده‌های منفی به تعداد داده‌های مثبت روش‌های مختلفی وجود دارد که رایج‌ترین آن، انتخاب به صورت تصادفی به تعداد داده‌های مثبت است. با انتخاب ۴۷۰ داده منفی از بین ۴۲۴۴ داده منفی، در نهایت مجموعه‌ای ۹۴۰ تایی با تعدادی مساوی از پپتیدهای ضدسرطانی و غیرضدسرطانی بدست آمد. نمایی کلی از مجموعه داده اصلی برای استفاده در ادامه پژوهش در شکل ۲-۱۰ مشاهده می‌شود.



شکل ۲-۹: نمایی کلی از مجموعه داده منفی cd-hit90 پس از انتخاب رندوم



شکل ۲-۱۰: نمایی کلی از مجموعه داده مثبت و منفی پس از اعمال پیش پردازش‌ها

۲-۳ مجموعه داده آموزش^{۳۴} و تست^{۳۵}

برای آموزش و ارزیابی مدل‌های یادگیری ماشین، مجموعه داده ۹۴۰ تایی از پیتیدهای بدست آمده در مرحله قبل به نسبت ۸۰-۲۰ طوری که ۸۰ درصد برای مجموعه داده آموزش و ۲۰ درصد برای مجموعه داده تست مستقل در نظر گرفته شد.

۲-۳-۱ مجموعه داده آموزش

برای یادگیری مدل‌های یادگیری ماشین و همچنین تنظیم پارامترهای^{۳۶} مدل‌ها، ۸۰ درصد از کل مجموعه داده مثبت یعنی ۳۷۶ داده مثبت و همچنین ۸۰ درصد از مجموعه داده منفی یعنی ۳۷۶ داده منفی به صورت کاملاً رندوم و بدون تکرار انتخاب شدند، که در مجموع برای هر بار آموزش و تنظیم پارامترها، ۷۵۲ داده کنار گذاشته شد.

۲-۳-۲ مجموعه داده تست مستقل^{۳۷}

برای ارزیابی عملکرد مدل‌های آموزش دیده ۲۰ درصد از ۹۴۰ پیتید موجود کنار گذاشته شد، طوری که این ۲۰ درصد در هیچ مرحله‌ای از آموزش و انتخاب پارامترهای تنظیمی استفاده نشدند. بعبارتی ۱۸۸ داده طوری که تعداد ۹۴ داده مثبت بصورت کاملاً رندوم و بدون تکرار از مجموعه داده‌های مثبت و همچنین تعداد ۹۴ داده منفی بصورت کاملاً رندوم و بدون تکرار از مجموعه داده‌های منفی انتخاب شدند.

³⁴ Training Dataset

³⁵ Test Set

³⁶ Hyper Parameter Tuning

³⁷ Independent Test Set

۲-۴ استخراج ویژگی^{۳۸}

با توجه به مطالعات پیشین در حوزه پیش‌بینی پپتیدهای ضدسرطانی، تصمیم به استخراج ویژگی از ساختار اولیه پپتیدها شد، با بررسی ابزارهای موجود برای استخراج ویژگی از توالی پپتیدها، کتابخانه پایتون iFeature نیازمندی‌های این پایان‌نامه را پوشش می‌دهد [۳۳]، در ادامه ۱۰ ویژگی استخراج شده مبتنی بر ساختار اول پپتیدها که در این پایان‌نامه استفاده شده است، به اختصار توضیح داده می‌شود [۳۴، ۳۵، ۳۷، ۳۸، ۴۲، ۴۳].

۲-۴-۱ ویژگی Pseudo-Amino Acid Composition (PseudoAAC)

فرض کنید $H_1^o(i)$ ، $H_2^o(i)$ و $M^o(i)$ برای $i = 1, 2, \dots, 20$ به ترتیب مقادیر آب‌گریزی، آب‌دوستی و وزن زنجیره جانبی برای ۲۰ اسید آمینه معمول باشند. مقدار آب‌گریزی برای هر اسید آمینه از رابطه زیر بدست می‌آید.

$$H_1(i) = \frac{H_1^o(i) - \frac{1}{20} \sum_{i=1}^{20} H_1^o(i)}{\sqrt{\frac{\sum_{i=1}^{20} [H_1^o(i) - \frac{1}{20} \sum_{i=1}^{20} H_1^o(i)]^2}{20}}} \quad \text{رابطه (۱-۲)}$$

$M(i)$ و $H_2(i)$ نیز به ترتیب برای ویژگی‌های آب‌دوستی و وزن زنجیره جانبی برای هر اسید آمینه مانند رابطه (۲-۱) به دست می‌آیند. سپس تابع correlation^{۳۹}ی به ازای هر جفت اسید آمینه R_i و R_j مقدار میانگین این ۳ ویژگی از رابطه (۲-۲) محاسبه می‌شود.

³⁸ Feature Extraction

$$\theta(R_i, R_j) = \frac{1}{3} \{ [H_1(i) - H_1(j)]^2 + [H_2(i) - H_2(j)]^2 + [M(i) - M(j)]^2 \} \quad \text{رابطه (۲-۱)}$$

سپس مجموعه عوامل مرتبط با ترتیب توالی^{۳۹} به صورت زیر تعریف می‌شود. N طول توالی و $\lambda < N$ است. در اینجا λ برابر با ۴ در نظر گرفته شده‌است.

$$\theta_1 = \frac{1}{N-1} \sum_{i=1}^{N-1} \theta(R_i, R_{i+1})$$

$$\theta_2 = \frac{1}{N-2} \sum_{i=1}^{N-2} \theta(R_i, R_{i+2})$$

رابطه (۳-۲)

⋮

$$\theta_\lambda = \frac{1}{N-\lambda} \sum_{i=1}^{N-\lambda} \theta(R_i, R_{i+\lambda})$$

در نهایت برداری به طول $20 + \lambda$ در اینجا به طول ۲۴ برای هر پروتئین از رابطه (۴-۲) بدست می‌آید که به آن ویژگی PseudoAAC گفته می‌شود. f_i فراوانی نرمال شده اسیدآمین به نام w و $w=0.5$ در نظر گرفته شده‌است [۳۳].

³⁹ Sequence Order-Correlated Factors

$$x_c = \frac{f_c}{\sum_{r=1}^{r_0} f_r + \omega \sum_{j=1}^{\lambda} \theta_j} \quad (1 < c < 20)$$

رابطه (۴-۲)

$$x_c = \frac{\omega \theta_{c-r_0}}{\sum_{r=1}^{r_0} f_r + \omega \sum_{j=1}^{\lambda} \theta_j} \quad (20 < c < 20 + \lambda)$$

۲-۴-۲ ویژگی APsedoAAC (Amphiphilic Pseudo-Amino Acid Composition)

تعریف مجموعه ویژگی‌های این ویژگی مشابه ویژگی PseudoAAC است. $H_1(i)$ و $H_2(i)$ ویژگی‌های آگریزی و آب‌دوستی از رابطه همبستگی زیر محاسبه می‌شود.

رابطه (۵-۲)

$$H_{i,j}^1 = H_1(i)H_1(j)$$

$$H_{i,j}^2 = H_2(i)H_2(j)$$

همچنین مجموعه عوامل مرتبط با ترتیب توالی به صورت زیر تعریف می‌شود.

رابطه (۶-۲)

$$\tau_1 = \frac{1}{N-1} \sum_{i=1}^{N-1} H_{i,i+1}^1$$

$$\tau_2 = \frac{1}{N-1} \sum_{i=1}^{N-1} H_{i,i+1}^2$$

$$\tau_3 = \frac{1}{N-2} \sum_{i=1}^{N-2} H_{i,i+2}^1$$

$$\tau_4 = \frac{1}{N-2} \sum_{i=1}^{N-2} H_{i,i+2}^2$$

...

$$\tau_{\lambda-1} = \frac{1}{N-\lambda} \sum_{i=1}^{N-\lambda} H_{i,i+\lambda}^1$$

$$\tau_{\lambda} = \frac{1}{N-\lambda} \sum_{i=1}^{N-\lambda} H_{i,i+\lambda}^{\gamma}$$

سپس مجموعه‌ای از توصیف‌گرها که به آن‌ها APseudoAAC گفته می‌شود به صورت زیر تعریف می‌شوند.

$$w=0.5$$

$$P_c = \frac{f_c}{\sum_{r=1}^{\gamma} f_r + w \sum_{j=1}^{\gamma\lambda} \tau_j}, (1 < c < \gamma)$$

$$P_c = \frac{\omega \tau_u}{\sum_{r=1}^{\gamma} f_r + w \sum_{j=1}^{\gamma\lambda} \tau_j}, (\gamma < u < \gamma + \gamma\lambda)$$

۲-۴-۳ ویژگی Composition of k-spaced Amino Acid Pairs (CKSAAP)

ویژگی CKSAAP، تعداد تکرار جفت اسیدآمینه که توسط k اسیدآمینه از هم فاصله γ دارند را محاسبه می‌کند. در این پایان نامه $k = 0, 1, 2, 3$ در نظر گرفته شده‌است. به همین دلیل بردار ویژگی به طول ۱۶۰۰ برای CKSAAP ساخته شد.

۲-۴-۴ ویژگی Composition/Transition/Distribution (CTD)

ویژگی CTD مخفف سه کلمه ترکیب (Composition)، توزیع (Distribution) و انتقال (Transition) است؛ این ویژگی الگوهای توزیع ویژگی‌های ساختاری یا فیزیکوشیمیایی یک اسیدآمینه در یک توالی

⁴⁰ Gap

پروتئین یا پپتید را نشان می‌دهد. برای محاسبه این ویژگی از ۱۳ ویژگی فیزیکوشیمیایی استفاده شده است؛ که شامل ویژگی‌های آبگریزی، قطبیت^{۴۱}، قطبش پذیری^{۴۲}، بار^{۴۳}، حجم نرمال شده وندروالس^{۴۴}، ساختارهای ثانویه^{۴۵} و قابلیت دسترسی به حلال^{۴۶} است. هر یک از این ویژگی‌ها بیست اسیدآمینه را بر اساس خواص فیزیکوشیمیایی به سه گروه تقسیم می‌کنند و سپس توالی پروتئینی بر اساس این که هر اسیدآمینه در کدام گروه قرار می‌گیرند بازنویسی می‌شوند. ویژگی‌های ترکیب، انتقال و توزیع برای توالی بازنویسی شده محاسبه می‌شود و هریک به ترتیب بردارهایی به طول ۳۹، ۳۹ و ۱۹۵ تولید می‌کنند. در نهایت بردار ویژگی CTD برای یک پروتئین از کنار هم قرار گرفتن این سه بردار ساخته می‌شود که به طول ۲۷۳ است.

⁴¹ Polarity

⁴² Polarizability

⁴³ Charge

⁴⁴ Normalized VanderWaals Volume

⁴⁵ Secondary Structures

⁴⁶ Solvent Accessibility

۵-۴-۲ ویژگی Dipeptide Deviation from Expected Mean (DDE)

ویژگی DDE توسط ۳ پارامتر ساخته می‌شود، ۱- ترکیب دوتایی پپتیدها $(D_c)^{47}$ ، میانگین نظری $(T_m)^{48}$ و واریانس نظری $(T_v)^{49}$. ۳ پارامترهای بیان شده و DDE در روابط زیر محاسبه می‌شوند.

$D_c(r,s)$ ، معیار ترکیب دوتایی برای دوپپتید 'rs':

$$D_c(r,s) = \frac{N_{rs}}{N-1}, r,s \in \{A,C,D,...Y\}$$

که، N_{rs} تعداد پپتیدهای دوتایی با اسیدآمینه r و s ، N طول پپتید هستند.

$T_m(r,s)$ ، میانگین نظری:

$$T_m(r,s) = \frac{C_r}{C_N} \times \frac{C_s}{C_N}$$

که، C_r تعداد کدون^{۵۰}های کدکننده اولین اسیدآمینه و C_s تعداد کدونهای کدکننده دومین اسیدآمینه در پپتید دوتایی 'rs' است. C_N تعداد کل کدونهای ممکن بغیر از ۳ کدون متوقف کننده^{۵۱} است.

$T_v(r,s)$ ، واریانس نظری:

$$T_v(r,s) = \frac{T_m(r,s)(1 - T_m(r,s))}{N-1}$$

و، در نهایت DDE(r,s) از رابطه زیر محاسبه می‌شود:

$$DDE(r,s) = \frac{D_c(r,s) - T_m(r,s)}{\sqrt{T_v(r,s)}}$$

⁴⁷ Dipeptide Composition

⁴⁸ Theoretical Mean

⁴⁹ Theoretical Variance

⁵⁰ Codon

⁵¹ Stop Codons

ویژگی DDE در نهایت بردار ویژگی به طول ۴۰۰ تولید می‌کند.

Moran correlation (Moran) ویژگی ۲-۴-۶

...
H CIDH920105
D Normalized average hydrophobicity scales (Cid et al., 1992)
R PMID:1518784
A Cid, H., Bunster, M., Canales, M. and Gazitua, F.
T Hydrophobicity and structural classes in proteins
J Protein Engineering 5, 373-375 (1992)

C	CIDH920103	0.973	CIDH920104	0.970	CIDH920102	0.969
	NISK860101	0.938	BASU050102	0.931	ZHOH040103	0.926
	ROBB790101	0.921	CIDH920101	0.921	MIYS850101	0.916
	BASU050103	0.914	PLIV810101	0.914	BIOV880101	0.912
	BASU050101	0.907	WERD780101	0.905	ZHOH040101	0.904
	RADA880108	0.898	FAUJ830101	0.893	MEEJ810101	0.892
	PONP930101	0.891	SWER830101	0.890	CORJ870102	0.890
	ROSM880104	0.886	BIOV880102	0.882	MANP780101	0.879
	ARGP820101	0.867	JOND750101	0.866	RADA880102	0.861
	CASG920101	0.859	GUOD860101	0.858	ROSG850102	0.858
	NOZY710101	0.857	PONP800101	0.856	NISK800101	0.854
	BLAS910101	0.852	CORJ870107	0.848	MEEJ810102	0.844
	PONP800108	0.843	ROSM880105	0.843	MEEJ800102	0.840
	TAKK010101	0.840	EISD860101	0.839	CORJ870104	0.838
	CORJ870103	0.838	SIMZ760101	0.837	PONP800102	0.831
	LIFS790101	0.828	LEVW760106	0.828	CORJ870101	0.827
	CORJ870106	0.826	CORJ870105	0.822	GOLD730101	0.820
	ZHOH040102	0.818	PONP800107	0.818	NADH010104	0.817
	PTIO830102	0.813	VENT840101	0.813	NADH010103	0.810
	PONP800103	0.807	MEIH800103	0.804	NADH010105	0.800
	WOEC730101	-0.800	KIDA850101	-0.803	PUNTO30101	-0.805
	KRIW790101	-0.816	FUKS010103	-0.821	PUNTO30102	-0.822
	MEIH800102	-0.826	RACS770102	-0.830	VINM940103	-0.832
	KARP850102	-0.839	CORJ870108	-0.843	FASG890101	-0.860
	PARS000101	-0.860	KARP850101	-0.866	BULH740101	-0.871
	GRAR740102	-0.884	VINM940101	-0.885	MIYS990103	-0.886
	RACS770101	-0.887	GUYH850102	-0.892	WOLS870101	-0.899
	MIYS990105	-0.901	O0BM770103	-0.904	MIYS990104	-0.908
	VINM940102	-0.910	MIYS990102	-0.915	MIYS990101	-0.916
	MEIH800101	-0.923	GUYH850103	-0.927	PARJ860101	-0.948

I	A/L	R/K	N/M	D/F	C/P	Q/S	E/T	G/W	H/Y	I/V
	0.02	-0.42	-0.77	-1.04	0.77	-1.10	-1.14	-0.80	0.26	1.81
	1.14	-0.41	1.00	1.35	-0.09	-0.97	-0.77	1.71	1.11	1.13

//

شکل ۲-۱۱: یک مثال مصور از ویژگی‌های فیزیکوشیمیایی در دیتابیس AAindex [۴۴]

توصیف کننده‌های^{۵۲} خودهمبستگی^{۵۳} براساس توزیع خواص^{۵۴} اسیدآمینها در توالی تعریف می‌شوند. خواص (ویژگی‌های) اسیدآمینها که در این توصیف کننده‌ها استفاده می‌شود، از دیتابیس AAindex استخراج می‌شود.

⁵² Descriptors

⁵³ Autocorrelation

⁵⁴ Properties

دیتابیس از مسیر <http://www.genome.jp/dbget/aaindex.html/> در دسترس است. ۸ شاخص

'CIDH920105', 'BHAR880101', 'CHAM820101', 'CHAM820102', 'CHOC760101',
'BIGC670101', 'CHAM810101', 'DAYM780201'

از دیتابیس AAindex استفاده شده‌اند.

تمام شاخص‌های اسیدآمین‌ها متمرکز^{۵۵} و استاندارد شده^{۵۶} توسط رابطه زیر هستند.

$$P_r = \frac{P_r - \bar{P}}{\sigma}$$

که، \bar{P} میانگین خواص ۲۰ اسیدآمین و σ انحراف معیار ۲۰ اسیدآمین است، که \bar{P} و σ خود توسط رابطه زیر قابل محاسبه هستند.

$$\bar{P} = \frac{\sum_{r=1}^{20} P_r}{20}, \sigma = \sqrt{\frac{1}{20} \sum_{r=1}^{20} (P_r - \bar{P})^2}$$

طول بردار ویژگی توصیف‌کننده‌های خود همبستگی $lag * n$ است، که n تعداد شاخص‌های انتخاب‌شده از پایگاه داده AAindex است. که در این پایان‌نامه $n=8$ و $lag=4$ و در نتیجه طول بردار ویژگی‌ها ۳۲ است. توصیف‌کننده خودهمبستگی Moran توسط رابطه زیر بدست می‌آید.

$$I(d) = \frac{\frac{1}{N-d} \sum_{i=1}^{N-d} (P_i - \bar{P}')(P_{i+d} - \bar{P}')}{\frac{1}{N} \sum_{i=1}^N (P_i - \bar{P}')^2}, d = 1, 2, 3, \dots, nlag$$

⁵⁵ Centralized

⁵⁶ Standardized

که، d lag خودهمبستگی و $nlag$ بیشترین مقدار برای lag است که در اینجا $nlag$ برابر با ۴ در نظر گرفته ایم.

P_i و P_{i+d} ویژگی‌های اسیدآمینها در موقعیت i و $i + d$ است. \bar{P} میانگین ویژگی P نیز از رابطه زیر بدست می‌آید.

$$\bar{P} = \frac{\sum_{i=1}^N P_i}{N}$$

۲-۴-۷ ویژگی Geary correlation (Geary)

برای یک پپتید یا پروتئین توصیف کننده‌های Geary از رابطه زیر محاسبه می‌شوند.

$$C(d) = \frac{\frac{1}{2(N-d)} \sum_{i=1}^{N-d} (P_i - P_{i+d})^2}{\frac{1}{N-1} \sum_{i=1}^N (P_i - \bar{P})^2}, d = 1, 2, \dots, nlag$$

که، d , P_i , P_{i+d} و $nlag$ تعاریف مشابه ویژگی Moran دارند.

۲-۴-۸ Normalized Moreau-Broto Autocorrelation (NMBroto) ویژگی

توصیف کننده‌های خودهمبستگی Moreau-Broto با رابطه زیر تعریف می‌شوند.

$$AC(d) = \sum_{i=1}^{N-d} P_i \times P_{i+d}, d = 1, 2, \dots, nlag$$

همچنین توصیف کننده‌های خودهمبستگی normalized Moreau-Broto با رابطه زیر تعریف می‌شوند.

$$ATS(d) = \frac{AC(d)}{N - d}, d = 1, 2, \dots, nlag$$

۹-۴-۲ ویژگی k -Spaced Conjoint Triad (KSCTriad)

توصیفگر KSCTriad مبتنی بر توصیفگر CTriad است که نه تنها اعداد سه واحد اسیدآمین به پیوسته را محاسبه می‌کند، بلکه واحدهای اسیدآمین به پیوسته را نیز در نظر می‌گیرد که با هر k باقیمانده^{۵۷} از هم جدا می‌شوند (حداکثر مقدار پیش‌فرض k در این پایان‌نامه ۱ تنظیم شده است). به عنوان مثال، $AxRxT$ یک سه‌گانه با $k=1$ است. بنابراین، ابعاد بردار ویژگی کدگذاری شده در KSCTriad، $343*(k+1) = 686$ است.

۱۰-۴-۲ ویژگی Quasi-sequence-order (QSOrder)

برای هر نوع اسیدآمین، یک توصیفگر QSOrder را می‌توان به صورت زیر تعریف کرد:

$$X_r = \frac{f_r}{\sum_{r=1}^{20} f_r + w \sum_{d=1}^{nlag} \tau_d}, r = 1, 2, \dots, 20$$

که در آن f_r وقوع نرمال شده اسیدآمین نوع r و w یک عامل وزنی است ($w=0.1$). اینها اولین ۲۰ توصیفگر QSOrder هستند. ۳۰ توصیفگر دیگر QSOrder به صورت زیر تعریف می‌شوند:

$$X_d = \frac{w\tau_d - 20}{\sum_{r=1}^{20} f_r + w \sum_{d=1}^{nlag} \tau_d}, d = 21, 22, \dots, 20 + nlag$$

۵-۲ پیش‌پردازش داده^{۵۸}

۴ مرحله تحت عنوان پیش‌پردازش بر روی داده‌های هر ۱۰ ویژگی استخراج شده اعمال شد.

حذف ستون‌های تکراری، از بین ستون‌هایی که مقادیر تکراری داشتند، همه بجز اولین ستون از بین تکراری‌ها حذف شدند.

حذف ستون با مقادیر یکسان، ستونی که برای تمام نمونه‌ها مقادیر یکسانی داشته باشد؛ حاوی اطلاعات مفید نیست به همین دلیل این ستون‌ها شناسایی و حذف شدند.

حذف نمونه‌های تکراری، نمونه‌هایی (سطرهایی) که برای تمام ویژگی‌ها مقدار یکسانی داشتند، همگی بجز اولین نمونه حذف شدند.

۶-۲ نرمالسازی داده^{۵۹}

از آنجایی که ممکن است در داده‌ها ویژگی‌های متفاوت محدوده‌های مختلفی داشته باشند، در این صورت ویژگی‌هایی که شامل دامنه اعداد بزرگتری هستند نسبت به بقیه ویژگی‌ها می‌توانند تأثیر بیشتری روی نتیجه پیش‌بینی داشته باشند که این باعث سو یافتن مدل یادگیری به سمت این ویژگی‌ها می‌شود در صورتی که ممکن است این ویژگی‌ها نسبت به ویژگی‌هایی که دامنه کوچکتری دارند اهمیت کمتری داشته باشند که این اتفاق خوبی نیست، به همین منظور تمام ویژگی‌ها باید نرمالسازی یا استانداردسازی شوند. در این مطالعه با استفاده از روش نرمال‌سازی min-max و با رابطه (۷-۲) هر ویژگی در مجموعه داده به بازه $[0, 1]$ نگاشت شدند.

⁵⁸ Data preprocessing

⁵⁹ Normalization

$$x' = \frac{x - x_{min}}{x_{max} - x_{min}} \quad \text{رابطه (۷-۲)}$$

هر یک از ویژگی‌های مجموعه داده تست مستقل و آموزش با رابطه (۷-۲) به اعداد جدیدی نگاشت شدند.

۷-۲ معیارهای ارزیابی

جهت سنجش کارایی الگوریتم‌های یادگیری ماشین نیاز به معیارهایی داریم و برای تعریف معیارهای ارزیابی ابتدا احتیاج به پارامترهایی داریم که بدانیم چه تعداد از نمونه‌ها به درستی پیش‌بینی شده‌اند. با توجه به این که این مطالعه در زمینه پیش‌بینی پتیدهای ضدسرطانی است، پارامترها و معیارهای ارزیابی بر این اساس تعریف خواهند شد.

مثبت صحیح^{۶۰}: یک پتید ضدسرطانی به درستی یک پتید ضدسرطانی اعلام شده‌است. تعداد این پیش‌بینی‌ها را با TP نشان می‌دهند.

مثبت کاذب^{۶۱}: یک پتید غیرضدسرطانی به اشتباه یک پتید ضدسرطانی اعلام شده‌است. تعداد این پیش‌بینی‌ها را با FP نشان می‌دهند.

منفی صحیح^{۶۲}: یک پتید غیرضدسرطانی به درستی یک پتید غیرضدسرطانی اعلام شده‌است. تعداد این پیش‌بینی‌ها را با TN نشان می‌دهند.

منفی کاذب^{۶۳}: یک پتید ضدسرطانی به اشتباه یک پتید غیرضدسرطانی اعلام شده‌است. تعداد این پیش‌بینی‌ها را با FN نشان می‌دهند.

⁶⁰ True positive

⁶¹ False positive

⁶² True negative

⁶³ False negative

در ادامه معیارهای ارزیابی براساس این پارامترها بیان می‌شود.

۱-۷-۲ حساسیت^{۶۴}

توانایی یک الگوریتم یادگیری ماشین برای تشخیص مقادیر مثبت صحیح (TP) را حساسیت یا یادآوری^{۶۵} گویند و از رابطه (۸-۲) بدست می‌آید. با توجه به این رابطه می‌توان دریافت که هر چه تعداد پتیدهای ضدسرطانی که به اشتباه پتید غیرضدسرطانی پیش‌بینی شده‌اند؛ کمتر باشد، یعنی منفی کاذب (FN) کمتر باشد، رابطه (۸-۲) به عدد یک نزدیک‌تر و حساسیت الگوریتم برای پیش‌بینی بالاتر است.

$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad \text{رابطه (۸-۲)}$$

۲-۷-۲ اختصاصیت^{۶۶}

توانایی یک الگوریتم یادگیری ماشین برای تشخیص مقادیر منفی صحیح (TN) را اختصاصیت گویند و از رابطه (۸-۲) بدست می‌آید. با توجه به این رابطه می‌توان دریافت که هر چه تعداد پتیدهای غیرضدسرطانی که به اشتباه پتید ضدسرطانی پیش‌بینی شده‌اند؛ کمتر باشد، یعنی مثبت کاذب (FP) کمتر باشد، رابطه (۹-۲) به عدد یک نزدیک‌تر اختصاصیت الگوریتم برای پیش‌بینی بالاتر است.

$$\text{Specificity} = \frac{TN}{TN + FP} \quad \text{رابطه (۹-۲)}$$

⁶⁴ Sensitivity

⁶⁵ Recall

⁶⁶ Specificity

۳-۷-۲ صحت^{۶۷}

یکی دیگر از معیارهای ارزیابی صحت است و از رابطه (۲-۸) بدست می‌آید. این معیار نشان می‌دهد چند درصد از داده‌ها به درستی طبقه‌بندی شده‌اند.

رابطه (۲-۱۰)

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

۴-۷-۲ دقت^{۶۸}

دقت نشان دهنده نسبت تعداد پتیده‌های ضدسرطانی که توسط الگوریتم یادگیری ماشین به درستی مثبت گزارش شده است به تعداد تمام پتیده‌های ضدسرطانی است؛ می‌باشد و با رابطه (۲-۱۱) محاسبه می‌شود.

$$Precision = \frac{TP}{TP + FP} \quad \text{رابطه (۲-۱۱)}$$

۸-۲ اعتبارسنجی متقابل k -لایه^{۶۹}

اعتبارسنجی متقابل k -لایه، یک روش برای ارزیابی مدل یادگیری ماشین است. در این روش مجموعه داده آموزش به k بخش افراز می‌شود. هر بار یک بخش به عنوان مجموعه داده آزمون در نظر گرفته می‌شود سپس مدل روی $k - 1$ بخش دیگر که به عنوان مجموعه آموزش انتخاب می‌شوند اجرا می‌شود. دقت مدل

⁶⁷ Accuracy

⁶⁸ Precision

⁶⁹ K-fold Cross Validation

روی بخشی که به عنوان مجموعه آزمون در نظر گرفته شده بود؛ محاسبه می‌شود. این روند k بار انجام می‌شود به طوری که تمام داده‌ها حداقل یکبار هم در مجموعه آزمون و هم در مجموعه آموزش شرکت کرده باشند و در نهایت میانگین دقت محاسبه شده در k مرحله به عنوان دقت مدل گزارش می‌شود. برای ارزیابی مدل‌های یادگیری در این پایان نامه مقدار k برابر ۱۰ در نظر گرفته شده است [۴۵].

۲-۹ انتخاب ویژگی^{۷۰}

انتخاب ویژگی، انتخاب یک زیرمجموعه از کل فضای مجموعه ویژگی‌ها است به طوری که ویژگی‌های انتخاب شده بیشترین ارتباط را با برجسب کلاس‌ها داشته باشند. انتخاب ویژگی‌های مناسب باعث می‌شود مدل‌های یادگیری ماشین کارایی بهتری داشته و ارتباط بین ویژگی‌ها و کلاس‌ها قابل درک‌تر باشد. به طور کلی روش‌های انتخاب ویژگی را می‌توان به دو گروه بسته‌بندی^{۷۱} و فیلتر^{۷۲} تقسیم کرد.

۲-۹-۱ روش بسته‌بندی

در این روش، مدلی با زیرمجموعه‌های مختلف از ویژگی‌ها توسط حذف یا اضافه کردن آن‌ها ساخته و آموزش داده می‌شود. در نهایت زیرمجموعه‌ای از ویژگی‌ها که بالاترین عملکرد را برای مدل داشتند، انتخاب می‌شوند. از روش‌های انتخاب ویژگی به روش بسته‌بندی می‌توان انتخاب متوالی رو به جلو^{۷۳} و انتخاب متوالی رو به

⁷⁰ Feature selection

⁷¹ Wrapper

⁷² Filter

⁷³ Sequential Forward Selection

عقب^{۷۴} را نام برد. در واقع این روش یک الگوریتم جستجو است که ویژگی‌ها به عنوان ورودی الگوریتم هستند و کارآیی مدل به عنوان خروجی؛ که باید بهینه شود.

در این پایان‌نامه از روش انتخاب متوالی رو به جلو توسط کتابخانه پایتون `mlxtend` استفاده شده است [۴۶].

انتخاب متوالی رو به جلو ۳ بار توسط پارامترهای زیر اجرا و ۳ دسته ویژگی توسط الگوریتم تولید شد.

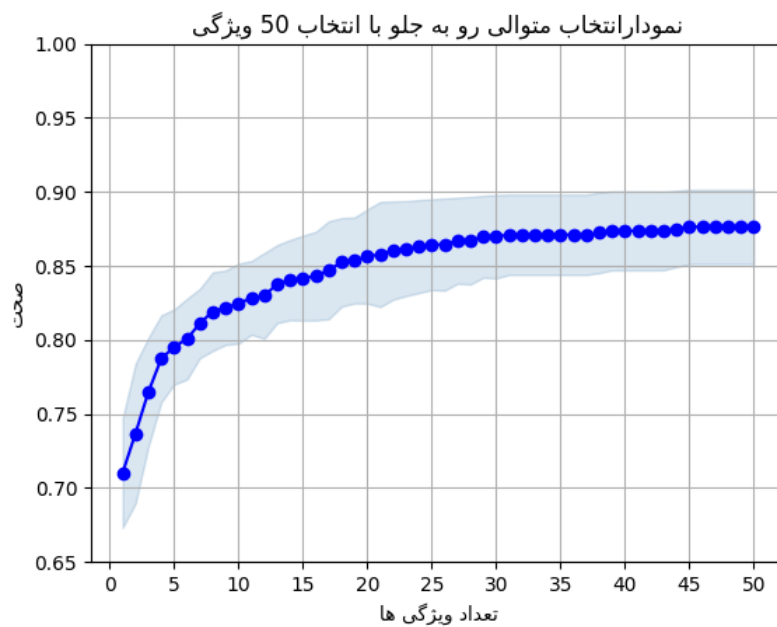
یکی از پارامترهای روش بسته‌بندی انتخاب یک مدل یادگیری ماشین برای انتخاب ویژگی‌ها است. در این پایان‌نامه از طبقه‌بند `Logestic Regression` برای این امر استفاده شد. معیار ارزیابی طبقه‌بند نیز صحت اعتبارسنجی ۱۰ لایه اعمال شد.

پارامتر بعدی در الگوریتم `SFS`^{۷۵}، برای تعداد ویژگی‌های انتخابی است، که در این جا ۳ مقدار ۵۰، ۱۰۰ و ۲۰۰ در نظر گرفته شد.

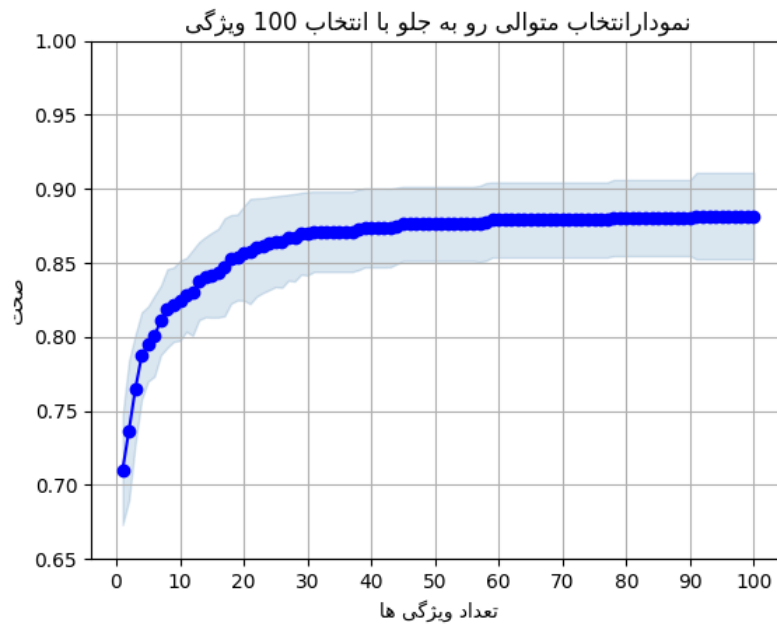
مراحل ۳ دسته ویژگی تولید شده که شامل ۵۰ ویژگی مهم، ۱۰۰ ویژگی مهم و ۲۰۰ ویژگی مهم هستند در تصاویر زیر مشاهده می‌شود. از این پس به ۵۰، ۱۰۰ و ۲۰۰ ویژگی انتخاب شده توسط الگوریتم `SFS`، با نام مجموعه داده `SFS50`، `SFS100` و `SFS200` اشاره خواهد شد.

⁷⁴ Sequential Backward Selection

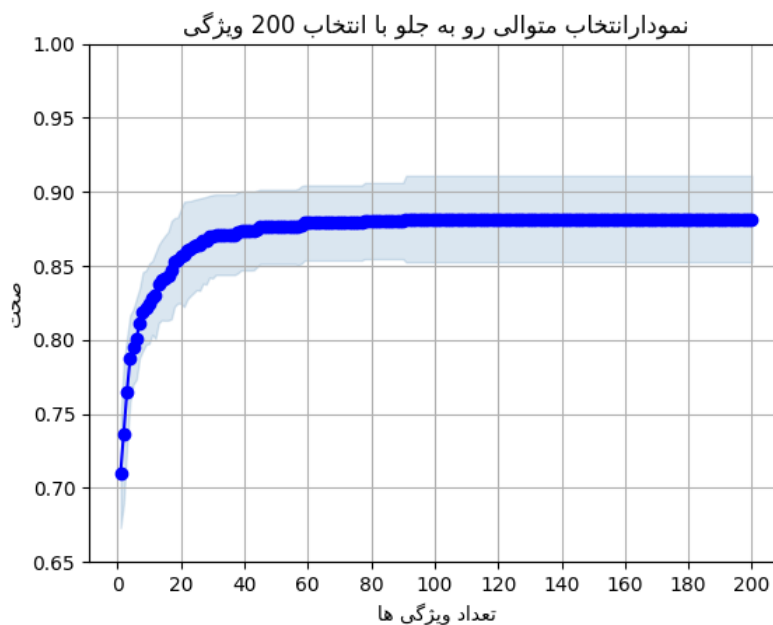
⁷⁵ Sequential Forward Selection



شکل ۲-۱۲: نمودار انتخاب متوالی رو به جلو تا ۵۰ ویژگی



شکل ۲-۱۳: نمودار انتخاب متوالی رو به جلو تا ۱۰۰ ویژگی



شکل ۲-۱۴: نمودار انتخاب متوالی رو به جلو تا ۲۰۰ ویژگی

در ادامه بر روی مجموعه داده SFS50، SFS100 و SFS200، مدل‌های یادگیری ماشین SVM و RF اجرا شد. نتایج این مدل‌ها در فصل ۳ آمده است.

۱۰-۲ کاهش ابعاد^{۷۶} مجموعه داده

تعداد زیاد ویژگی‌ها باعث تحمیل هزینه محاسباتی بالا به طبقه‌بندها می‌شود و به مشقت تعداد ابعاد^{۷۷} معروف است. برای کاهش ابعاد می‌توان از روش تحلیل مؤلفه‌های اصلی^{۷۸} استفاده کرد.

⁷⁶ Dimension Reduction

⁷⁷ Curse of Dimensionality

⁷⁸ Principal Component Analysis

به صورت ریاضی تحلیل مؤلفه‌های اصلی (PCA) یک تبدیل خطی متعامد است که داده‌ها را به دستگاه مختصات جدید می‌برد. به صورتی که اولین بزرگترین واریانس داده بر روی اولین محور مختصات، دومین بزرگترین واریانس بر روی دومین محور مختصات و ... قرار می‌گیرد [۴۷].

در این پایان‌نامه، این روش برای تمام ویژگی‌هایی که تا به حال بحث شده‌است با انتخاب ۱۰ مؤلفه اصلی اول اعمال شد، در ادامه بر روی این ۱۰ مؤلفه (۱۰ ویژگی) مدل‌های یادگیری ماشین SVM و RF اجرا شد. نتایج این مدل‌ها در فصل ۳ آمده‌است.

۱۱-۲ مدل‌های یادگیری ماشین اجرا شده روی مجموعه داده‌ها

همانطور که در بخش ۲-۲ بیان شد، بر روی مجموعه داده cdhit90 پیش‌پردازش اولیه انجام شد؛ سپس در بخش استخراج ویژگی (بخش ۲-۴) به ۱۰ ویژگی اشاره شد.

این ۱۰ ویژگی و همچنین ترکیبی از این ۱۰ ویژگی که در ادامه اشاره خواهد شد، از مجموعه داده پیش-پردازش شده بخش ۲-۲، استخراج شدند. سپس ۱۰ مؤلفه اول هر یک از این مجموعه داده‌های بدست‌آمده توسط الگوریتم PCA از کتابخانه iFeature استخراج شد.

در نهایت بر روی هر کدام از مجموعه داده‌های بدست‌آمده توسط PCA و همچنین بر روی ۳ مجموعه داده SFS50، SFS100 و SFS200 (رجوع به بخش ۲-۹-۱) ۲ مدل یادگیری ماشین SVM و RF با پارامترهای تنظیمی اشاره شده در قسمت بعدی ساخته شدند. ارزیابی مدل‌ها و یافتن بهترین پارامتر تنظیمی از بین پارامترهای تعریف شده توسط اعتبارسنجی متقابل ۱۰ لایه انجام شد. در تمام مراحل از

تعریف پارامترهای تنظیمی، تنظیم پارامترهای تنظیمی، آموزش و ارزیابی مدل‌ها؛ از کتابخانه scikit-learn پایتون استفاده شد.

۱۱-۲ طبقه‌بندی جنگل تصادفی

برای ساخت هر جنگل تصادفی ۴۳۲۰، ۴۳۲۰ = ۱۰ * ۲ * ۱۲ * ۳ * ۳ * ۲ پارامتر تنظیمی، توسط قطعه کد زیر ایجاد شد.

الگوریتم انتخاب تصادفی پارامترها، از بین ۴۳۲۰ پارامتر تنظیمی، یک جنگل تصادفی بصورت کاملاً تصادفی می‌سازد. برای یافتن بهترین پارامتر تنظیمی، ۱۰۰ جنگل تصادفی تصادفی^{۷۹}، توسط پارامترهای تصادفی ایجاد شده، ساخته شد. [۴۸، ۴۹]

قطعه کد زیر، فضای پارامترهای تنظیمی برای ساخت ۱۰۰ جنگل تصادفی تصادفی را نشان می‌دهد.

```
# Number of trees in random forest
n_estimators = [int(x) for x in np.linspace(start=100, stop=1000, num=10)]
# Number of features to consider at every split
max_features = ['auto', 'sqrt']
# Maximum number of levels in tree
max_depth = [int(x) for x in np.linspace(10, 110, num=11)]
max_depth.append(None)
# Minimum number of samples required to split a node
min_samples_split = [2, 5, 10]
# Minimum number of samples required at each leaf node
min_samples_leaf = [1, 2, 4]
# Method of selecting samples for training each tree
bootstrap = [True, False]
# On each iteration, the algorithm will choose a difference combination of the features.
# Altogether, there are 2 * 12 * 2 * 3 * 3 * 10 = 4320 settings!
```

⁷⁹ Random Random Forest

۲-۱۱-۲ طبقه‌بندی ماشین بردار پشتیبان

برای ساخت هر ماشین بردار پشتیبان ۵۲، ۵۲ = ۱۳ * ۲ * ۲ پارامتر تنظیمی، توسط قطعه کد زیر ایجاد شد.

```
kernel = ['rbf', 'linear']
gamma = ['scale', 'auto']
c_range = [x for x in np.logspace(-9, 3, 13)]
```

kernel، وظیفه اصلی هسته^{۸۰} این است که مجموعه داده ورودی داده شده را به فرم مورد نیاز تبدیل کند. انواع مختلفی از توابع مانند تابع پایه خطی^{۸۱}، چند جمله‌ای^{۸۲} و شعاعی (RBF)^{۸۳} وجود دارد. تابع پایه چند جمله‌ای و RBF برای ابر صفحه غیر خطی مفید هستند. هسته‌های چند جمله‌ای و RBF خط جداسازی را در بعد بالاتر محاسبه می‌کنند. در اینجا ۲ پارامتر هسته خطی و هسته شعاعی برای پارامترهای هسته در نظر گرفته شد.

c_range، پارامتر هزینه^{۸۴} یا مجازات^{۸۵}، نشان دهنده طبقه‌بندی اشتباه^{۸۶} یا عبارت خطا^{۸۷} است. پارامتر c به بهینه‌سازی SVM می‌گوید که چقدر خطا قابل تحمل است. به این ترتیب می‌تواند trade-off بین مرز تصمیم^{۸۸} و اصطلاح طبقه‌بندی اشتباه را کنترل کند. مقدار کوچکتر c یک ابر صفحه با حاشیه کوچک و مقدار بزرگتر c یک ابر صفحه با حاشیه بزرگتر ایجاد می‌کند [۵۰، ۵۱].

gamma، مقدار کمتر پارامتر گاما با مجموعه داده آموزشی متناسب است، در حالی که مقدار بالاتر گاما دقیقاً با مجموعه داده آموزشی مطابقت دارد که باعث برآزش بیش از حد^{۸۹} می‌شود. به عبارت دیگر، می‌توان

⁸⁰ Kernel

⁸¹ Linear Basis Function

⁸² Polynomial Basis Function

⁸³ Radial Basis Function (rbf)

⁸⁴ Cost

⁸⁵ Penalty

⁸⁶ Missclassification

⁸⁷ Error Term

⁸⁸ Decision Boundary

⁸⁹ Over-Fitting

گفت مقدار کم گاما فقط نقاط نزدیک را در محاسبه خط جداسازی در نظر می‌گیرد، در حالی که مقدار بالای گاما تمام نقاط داده را در محاسبه خط جداسازی در نظر می‌گیرد.

برای انتخاب بهترین پارامترتنظیمی، ۵۲ مدل ماشین بردار پشتیبان، به تعداد پارامترهای تنظیمی تعریف شده، ساخته شد.

۳ فصل سوم: نتایج و بحث

۳-۱ نتایج طبقه‌بندیها

برای کاهش ابعاد داده‌های ۱۰ ویژگی استخراج شده و پیش‌پردازش شده، ۲ روش در نظر گرفته شد.

۱- استفاده از PCA و استخراج ۱۰ مؤلفه اول از داده‌ها آن‌ها

۲- استفاده از الگوریتم انتخاب ویژگی SFS، که نتیجه آن ۳ دسته مجموعه داده SFS50،

SFS100 و SFS200 شد.

سپس هر مجموعه داده بدست آمده به این ۲ روش، به نسبت ۸۰-۲۰ تقسیم شدند، طوریکه ۸۰

درصد داده‌ها برای آموزش و تنظیم پارامترهای تنظیمی و ۲۰ درصد داده‌ها برای تست مستقل در

نظر گرفته شدند. سپس ۲ طبقه‌بند جنگل تصادفی و ماشین بردار پشتیبان بر روی ۸۰ درصد داده‌های آموزش توسط روش اعتباری‌سنجی متقابل ۱۰ لایه، با پارامترهای تنظیمی تعریف شده در بخش ۲-۱۱ و ۲-۱۱-۲ تنظیم پارامتر^{۹۰} شدند.

مدلی که بالاترین عملکرد (صحت) را داشته به معنی این است که پارامترهای تنظیمی آن مدل، بر روی این مجموعه داده بهترین پارامترها هستند، این مدل‌ها به عنوان بهترین مدل انتخاب شدند. برای بار دوم، اینبار بهترین مدل‌ها با تمام ۸۰ درصد داده آموزش دیدند، سپس توسط مجموعه تست مستقل عملکرد مدل‌ها ارزیابی شد.

نتایج بهترین مدل‌ها توسط اعتبارسنجی متقابل ۱۰ لایه در مرحله تنظیم پارامترهای تنظیمی و نتایج مدل‌ها توسط تست مستقل در ادامه گزارش می‌شوند.

۳-۲ نتایج بدست آمده بر روی داده‌های ۱۰ مؤلفه اول PCA

در ادامه نتایج ۲ طبقه‌بند جنگل تصادفی و ماشین بردار پشتیبان بر روی ۳ سری مجموعه داده گزارش می‌شود.

سری ۱- داده‌های ۱۰ مؤلفه اول PCA که از تک‌ویژگی‌ها استخراج شد.

در این سری، بر روی ۱۰ ویژگی استخراج شده، به تفکیک ویژگی‌ها PCA و سپس RF و SVM اجرا شد.

سری ۲- داده‌های ۱۰ مؤلفه اول PCA که از ترکیب دوتایی ویژگی‌ها استخراج شد.

⁹⁰ Parameter Tuning

در این سری، ویژگی‌هایی که در سری ۱، در تست مستقل، صحت آن‌ها بالاتر از ۰.۷۳ گزارش شده، انتخاب شدند؛ که شامل ۷ دسته ویژگی شد. سپس تمام ترکیب‌های ۲ تایی این ۷ دسته ویژگی ساخته شد و توسط PCA، ۱۰ مؤلفه اول از آن‌ها استخراج شد. سپس بر روی این ۱۰ مؤلفه همانند روش قبل، ۲ مدل یادگیری ماشین RF و SVM اجرا شد.

سری ۳- داده‌های ۱۰ مؤلفه اول PCA که از ترکیب تمام ویژگی‌ها استخراج شد.

در این سری نیز، همانند سری ۲، ویژگی‌هایی که در سری ۱، در تست مستقل، صحت آن‌ها بالاتر از ۰.۷۳ گزارش شده بود، انتخاب شدند؛ که شامل ۷ دسته ویژگی شد. سپس تمام این ۷ ویژگی‌ها باهم ترکیب شدند و توسط PCA، ۱۰ مؤلفه اول از آن‌ها استخراج شد. سپس بر روی این ۱۰ مؤلفه همانند روش قبل، ۲ مدل یادگیری ماشین RF و SVM اجرا شد.

در ادامه نتایج بدست آمده از ۲ طبقه‌بند جنگل تصادفی و ماشین بردار پشتیبان بر روی این ۳ سری داده گزارش می‌شود.

۱-۲-۳ نتایج مربوط به طبقه‌بند جنگل تصادفی

۱-۱-۲-۳ نتایج مربوط به مجموعه داده‌های سری ۱

نتایج بهترین مدل جنگل تصادفی بر روی داده‌های سری ۱، داده‌های ۱۰ مؤلفه اول PCA که از تک‌ویژگی‌ها استخراج شد، در مرحله تنظیم پارامترهای تنظیمی و تست مستقل در جدول زیر گزارش می‌شود.

جدول ۳-۱: نتایج بهترین مدل جنگل تصادفی در هر تک‌ویژگی

ویژگی	صحت بهترین مدل در تنظیم پارامتر تنظیمی	معیار roc_auc بهترین مدل در تنظیم پارامتر تنظیمی	صحت (تست مستقل)	دقت (تست مستقل)	حساسیت (تست مستقل)	معیار f1 (تست مستقل)
APseudoAAC	۷۹.۰	۸۷.۰	۷۹.۰	۷۹.۰	۷۹.۰	۷۹.۰
CKSAAP	۸۰.۰	۸۹.۰	۷۹.۰	۷۸.۰	۷۸.۰	۷۸.۰
CTD	۷۶.۰	۸۴.۰	۷۴.۰	۷۴.۰	۷۴.۰	۷۴.۰
DDE	۷۷.۰	۸۵.۰	۷۴.۰	۷۴.۰	۷۴.۰	۷۴.۰
Geary	۷۱.۰	۷۷.۰	۷۱.۰	۷۱.۰	۷۲.۰	۷۱.۰
KSCTriad	۷۱.۰	۸۴.۰	۷۷.۰	۷۶.۰	۷۶.۰	۷۶.۰
Moran	۷۰.۰	۷۵.۰	۶۶.۰	۶۷.۰	۶۷.۰	۶۶.۰
NMBroto	۷۱.۰	۷۸.۰	۶۶.۰	۶۶.۰	۶۶.۰	۶۶.۰
PseudoAAC	۷۹.۰	۸۷.۰	۷۹.۰	۷۸.۰	۷۸.۰	۷۸.۰
QSOrder	۷۹.۰	۸۶.۰	۸۱.۰	۸۱.۰	۸۱.۰	۸۱.۰

۳-۲-۱-۲ نتایج مربوط به مجموعه داده‌های سری ۲

با توجه به جدول ۳-۱، ویژگی‌هایی که صحت تست مستقل آن‌ها بالاتر از ۰.۷۳ گزارش شده، انتخاب شدند؛ که شامل ۷ ویژگی APseudoAAC، CKSAAP، CTD، DDE، KSCTriad، PseudoAAC و QSOrder شد. سپس تمام ترکیب‌های ۲تایی این ۷ دسته ویژگی ساخته شد و توسط PCA، ۱۰ مؤلفه اول از آن‌ها استخراج شد. نتایج بهترین مدل‌های جنگل تصادفی بر روی این سری از داده‌ها، در مرحله تنظیم پارامترهای تنظیمی و تست مستقل در جدول زیر گزارش می‌شود.

جدول ۳-۲: نتایج بهترین مدل جنگل تصادفی در هر ترکیب دوتایی ویژگی‌ها

ویژگی	صحت (بهترین مدل در تنظیم پارامتر تنظیمی)	معیار roc_auc (بهترین مدل در تنظیم پارامتر تنظیمی)	صحت (تست مستقل)	دقت (تست مستقل)	حساسیت (تست مستقل)	معیار f1 (تست مستقل)
QOrder+CKSAAP	۸۲.۰	۸۹.۰	۸۰.۰	۸۰.۰	۸۰.۰	۸۰.۰
QOrder+CTD	۷۶.۰	۸۴.۰	۷۶.۰	۷۶.۰	۷۵.۰	۷۵.۰
QOrder+DDE	۸۳.۰	۸۹.۰	۸۰.۰	۷۹.۰	۸۰.۰	۷۹.۰
QOrder+KSCTriad	۷۹.۰	۸۶.۰	۷۷.۰	۷۷.۰	۷۶.۰	۷۷.۰
QOrder+PseudoAAC	۸۱.۰	۸۸.۰	۸۱.۰	۸۱.۰	۸۱.۰	۸۱.۰
QOrder+APseudoAAC	۸۱.۰	۸۸.۰	۸۱.۰	۸۱.۰	۸۱.۰	۸۱.۰
PseudoAAC+APseudoAAC	۷۹.۰	۸۷.۰	۷۹.۰	۷۹.۰	۷۹.۰	۷۹.۰
PseudoAAC+CKSAAP	۸۲.۰	۹۰.۰	۷۸.۰	۷۷.۰	۷۸.۰	۷۷.۰
PseudoAAC+CTD	۷۵.۰	۸۴.۰	۷۶.۰	۷۵.۰	۷۵.۰	۷۵.۰
PseudoAAC+DDE	۸۱.۰	۸۹.۰	۸۰.۰	۸۰.۰	۸۰.۰	۸۰.۰
PseudoAAC+KSCTriad	۷۸.۰	۸۶.۰	۷۶.۰	۷۵.۰	۷۵.۰	۷۵.۰
CKSAAP+APseudoAAC	۸۱.۰	۹۰.۰	۷۸.۰	۷۸.۰	۷۸.۰	۷۸.۰
CKSAAP+CTD	۷۸.۰	۸۶.۰	۷۷.۰	۷۷.۰	۷۶.۰	۷۷.۰
CKSAAP+DDE	۸۰.۰	۸۹.۰	۷۵.۰	۷۵.۰	۷۴.۰	۷۴.۰
CKSAAP+KSCTriad	۸۰.۰	۸۸.۰	۷۸.۰	۷۷.۰	۷۷.۰	۷۷.۰
CTD+APseudoAAC	۷۵.۰	۸۴.۰	۷۷.۰	۷۷.۰	۷۷.۰	۷۷.۰
CTD+DDE	۷۶.۰	۸۵.۰	۷۸.۰	۷۸.۰	۷۸.۰	۷۸.۰
CTD+KSCTriad	۷۸.۰	۸۶.۰	۸۰.۰	۸۰.۰	۷۹.۰	۸۰.۰
DDE+APseudoAAC	۸۱.۰	۸۹.۰	۸۰.۰	۸۰.۰	۸۰.۰	۸۰.۰
DDE+KSCTriad	۷۷.۰	۸۶.۰	۷۵.۰	۷۵.۰	۷۴.۰	۷۴.۰
APseudoAAC+KSCTriad	۷۹.۰	۸۶.۰	۷۷.۰	۷۶.۰	۷۶.۰	۷۶.۰

۳-۱-۲-۳ نتایج مربوط به ترکیب تمام ویژگی‌ها

با توجه به جدول ۱-۳، تمام ۷ ویژگی‌هایی که صحت تست مستقل آن‌ها بالاتر از ۰.۷۳ گزارش شده، باهم ترکیب شدند؛ سپس توسط PCA، ۱۰ مؤلفه اول از آن‌ها استخراج شد. نتایج بهترین مدل جنگل تصادفی بر روی این سری از داده، در مرحله تنظیم پارامترهای تنظیمی و تست مستقل در جدول زیر گزارش می‌شود.

جدول ۳-۳: نتایج بهترین مدل جنگل تصادفی در ترکیب تمام ۷ ویژگی‌ها

ویژگی	صحت (بهترین مدل در تنظیم پارامتر تنظیمی)	معیار roc_auc (بهترین مدل در تنظیم پارامتر تنظیمی)	صحت (تست مستقل)	دقت (تست مستقل)	حساسیت (تست مستقل)	معیار f1 (تست مستقل)
ترکیب همه ۷ ویژگی‌ها	۸۰.۰	۸۹.۰	۸۲.۰	۸۲.۰	۸۲.۰	۸۲.۰

۳-۲-۲ نتایج مربوط به طبقه‌بند ماشین بردار پشتیبان

۱-۲-۲-۳ نتایج مربوط به مجموعه داده‌های سری ۱

نتایج بهترین مدل‌های ماشین بردار پشتیبان بر روی داده‌های سری ۱، داده‌های ۱۰ مؤلفه اول PCA که از تک‌ویژگی‌ها استخراج شد، در مرحله تنظیم پارامترهای تنظیمی و تست مستقل در جدول زیر گزارش می‌شود.

جدول ۳-۴: نتایج بهترین مدل ماشین بردار پشتیبان در هر تک‌ویژگی

ویژگی	صحت (بهترین مدل در تنظیم پارامتر تنظیمی)	صحت (تست مستقل)	دقت (تست مستقل)	حساسیت (تست مستقل)	معیار f1 (تست مستقل)
APseudoAAC	۷۹.۰	۷۷.۰	۷۷.۰	۷۶.۰	۷۶.۰

۷۷.۰	۷۷.۰	۷۷.۰	۷۸.۰	۷۹.۰	CKSAAP
۷۷.۰	۷۷.۰	۷۷.۰	۷۸.۰	۷۷.۰	CTD
۷۳.۰	۷۳.۰	۷۳.۰	۷۳.۰	۷۷.۰	DDE
۶۵.۰	۶۵.۰	۶۵.۰	۶۵.۰	۶۹.۰	Geary
۷۲.۰	۷۳.۰	۷۲.۰	۷۳.۰	۷۷.۰	KSCTriad
۶۳.۰	۶۳.۰	۶۳.۰	۶۴.۰	۶۹.۰	Moran
۷۵.۰	۷۵.۰	۷۵.۰	۷۵.۰	۷۲.۰	NMBroto
۷۶.۰	۷۶.۰	۷۶.۰	۷۷.۰	۸۰.۰	PseudoAAC
۸۰.۰	۸۰.۰	۸۰.۰	۸۰.۰	۷۹.۰	QOrder

۲-۲-۲-۳ نتایج مربوط به مجموعه داده‌های سری ۲

با توجه به جدول ۳-۴، ویژگی‌هایی که صحت تست مستقل آن‌ها بالاتر از ۰.۷۳ گزارش شده، انتخاب شدند؛ که شامل ۸ ویژگی APseudoAAC، CKSAAP، CTD، DDE، KSCTriad، NMBroto، PseudoAAC و QOrder شد. سپس تمام ترکیب‌های ۲ تایی این ۸ دسته ویژگی ساخته شد و توسط PCA، ۱۰ مؤلفه اول از آن‌ها استخراج شد. نتایج بهترین مدل‌های ماشین بردار پشتیبان بر روی این سری از داده‌ها، در مرحله تنظیم پارامترهای تنظیمی و تست مستقل در جدول زیر گزارش می‌شود.

جدول ۳-۵: نتایج بهترین مدل جنگل تصادفی در هر ترکیب دوتایی ویژگی‌ها

ویژگی	صحت (بهترین مدل در تنظیم پارامتر تنظیمی)	صحت (تست مستقل)	دقت (تست مستقل)	حساسیت (تست مستقل)	معیار fl (تست مستقل)
QOrder+CKSAAP	۸۲.۰	۷۸.۰	۷۸.۰	۷۸.۰	۷۸.۰
QOrder+CTD	۷۷.۰	۷۷.۰	۷۷.۰	۷۶.۰	۷۷.۰
QOrder+DDE	۸۱.۰	۷۹.۰	۷۹.۰	۷۹.۰	۷۹.۰
QOrder+KSCTriad	۷۸.۰	۷۱.۰	۷۱.۰	۷۱.۰	۷۰.۰
QOrder+PseudoAAC	۸۱.۰	۷۹.۰	۷۹.۰	۷۹.۰	۷۹.۰

۸۱.۰	۸۱.۰	۸۱.۰	۸۱.۰	۸۱.۰	QOrder+APseudoAAC
۷۶.۰	۷۶.۰	۷۷.۰	۷۷.۰	۷۹.۰	PseudoAAC+APseudoAAC
۸۰.۰	۸۰.۰	۸۰.۰	۸۰.۰	۸۱.۰	PseudoAAC+CKSAAP
۷۵.۰	۷۵.۰	۷۵.۰	۷۶.۰	۷۶.۰	PseudoAAC+CTD
۷۳.۰	۷۴.۰	۷۳.۰	۷۳.۰	۸۰.۰	PseudoAAC+DDE
۷۱.۰	۷۱.۰	۷۱.۰	۷۱.۰	۷۸.۰	PseudoAAC+KSCTriad
۷۹.۰	۷۹.۰	۷۹.۰	۸۰.۰	۸۱.۰	CKSAAP+APseudoAAC
۷۴.۰	۷۴.۰	۷۵.۰	۷۵.۰	۷۸.۰	CKSAAP+CTD
۷۷.۰	۷۶.۰	۷۸.۰	۷۸.۰	۸۰.۰	CKSAAP+DDE
۷۵.۰	۷۵.۰	۷۵.۰	۷۵.۰	۷۹.۰	CKSAAP+KSCTriad
۷۵.۰	۷۵.۰	۷۵.۰	۷۶.۰	۷۷.۰	CTD+APseudoAAC
۷۷.۰	۷۷.۰	۷۷.۰	۷۸.۰	۷۶.۰	CTD+DDE
۷۲.۰	۷۲.۰	۷۲.۰	۷۲.۰	۷۷.۰	CTD+KSCTriad
۷۴.۰	۷۴.۰	۷۵.۰	۷۵.۰	۸۱.۰	DDE+APseudoAAC
۷۴.۰	۷۴.۰	۷۳.۰	۷۴.۰	۷۸.۰	DDE+KSCTriad
۷۳.۰	۷۴.۰	۷۳.۰	۷۴.۰	۷۸.۰	APseudoAAC+KSCTriad
۷۷.۰	۷۷.۰	۷۷.۰	۷۸.۰	۸۲.۰	NMBroto+APseudoAAC
۷۸.۰	۷۸.۰	۷۹.۰	۷۹.۰	۸۰.۰	NMBroto+CKSAAP
۷۴.۰	۷۴.۰	۷۵.۰	۷۵.۰	۷۶.۰	NMBroto+CTD
۷۶.۰	۷۵.۰	۷۶.۰	۷۷.۰	۸۰.۰	NMBroto+DDE
۷۳.۰	۷۳.۰	۷۳.۰	۷۳.۰	۷۸.۰	NMBroto+KSCTriad
۷۸.۰	۷۸.۰	۷۸.۰	۷۹.۰	۸۲.۰	NMBroto+PseudoAAC
۷۹.۰	۷۹.۰	۷۹.۰	۷۹.۰	۸۰.۰	NMBroto+QOrder

۳-۲-۲-۳ نتایج مربوط به ترکیب تمام ویژگی‌ها

با توجه به جدول ۳-۴، تمام ۸ ویژگی‌هایی که صحت تست مستقل آن‌ها بالاتر از ۰.۷۳ گزارش شده، باهم

ترکیب شدند؛ سپس توسط PCA، ۱۰ مؤلفه اول از آن‌ها استخراج شد. نتایج بهترین مدل‌های ماشین‌بردار

پشتیبان بر روی این سری از داده‌ها، در مرحله تنظیم پارامترهای تنظیمی و تست مستقل در جدول زیر گزارش می‌شود.

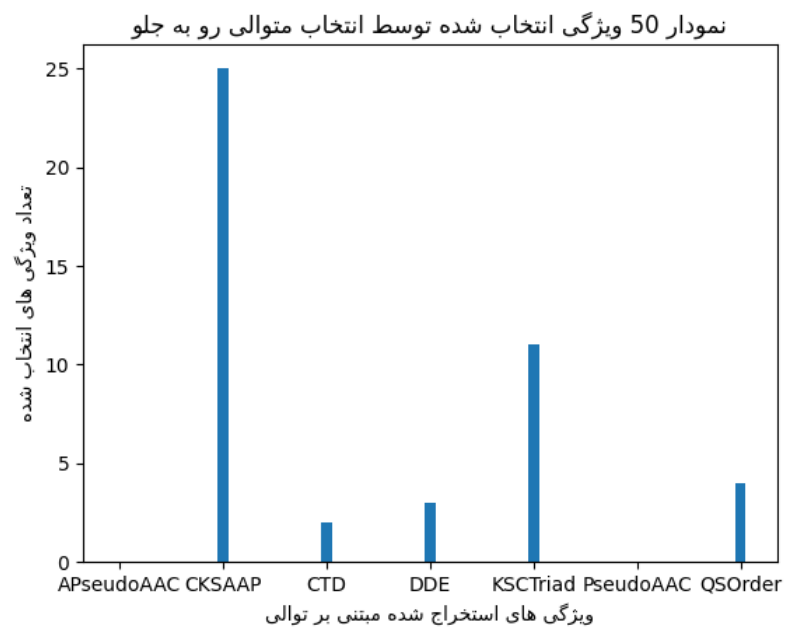
جدول ۳-۶: نتایج بهترین مدل ماشین بردار پشتیبان در ترکیب تمام ۸ ویژگی‌ها

ویژگی	صحت (بهترین مدل در تنظیم پارامتر تنظیمی)	صحت (تست مستقل)	دقت (تست مستقل)	حساسیت (تست مستقل)	معیار f1 (تست مستقل)
ترکیب همه ۸ ویژگی‌ها	۷۹.۰	۷۹.۰	۷۹.۰	۷۹.۰	۷۹.۰

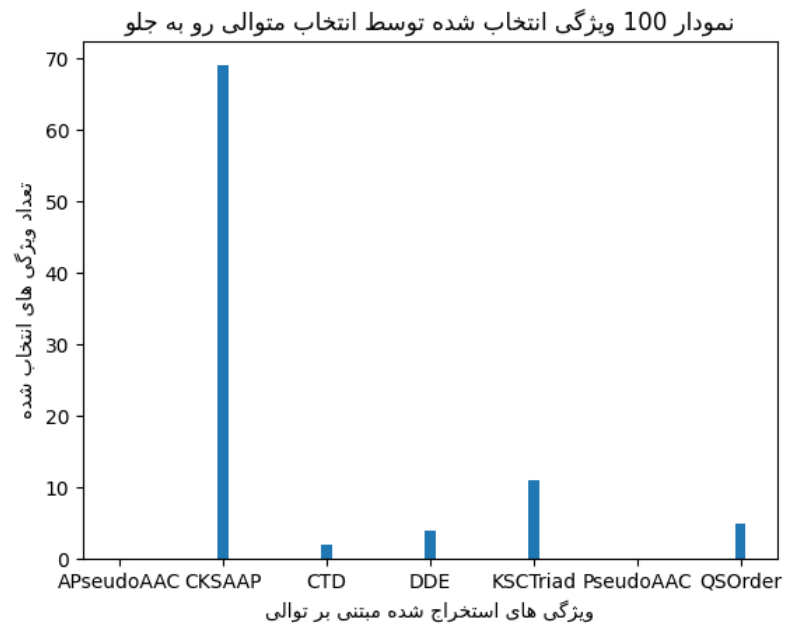
۳-۳ نتایج ویژگی‌های انتخاب شده توسط الگوریتم انتخاب متوالی رو به جلو

با کنار هم قراردادن تمام ۷ ویژگی مشخص شده در جدول ۳-۱ یک فضای ویژگی به ابعاد (۳۰۱۴، ۹۴۰) ایجاد شد. (۳۰۱۴ طول بردار ویژگی و ۹۴۰ تعداد نمونه پیتیدها)

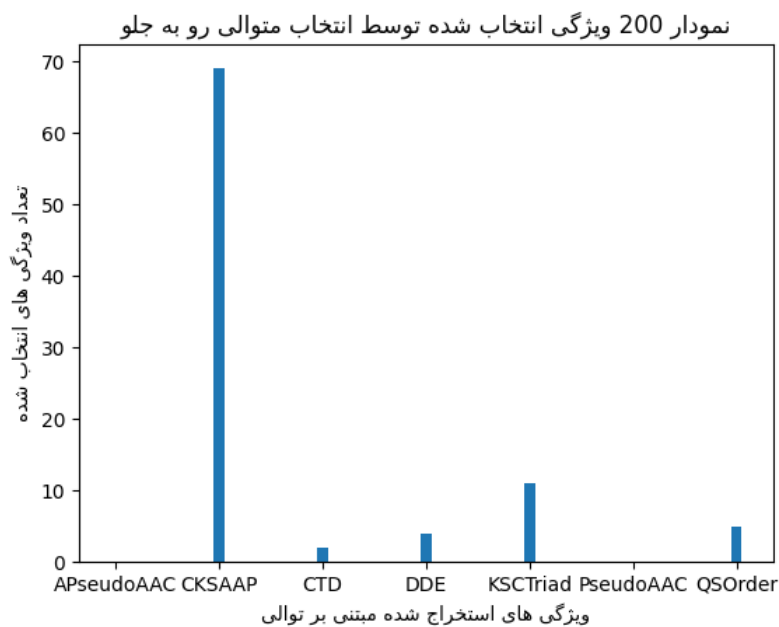
فضای ویژگی جدید به الگوریتم SFS داده شد و ۳ دسته ویژگی SFS50، SFS100 و SFS200 ساخته شد. برای درک این موضوع که چه تعداد از ۷ ویژگی در ویژگی‌های SFS50، SFS100 و SFS200 مشارکت داشته‌اند، نمودارهای زیر رسم شد. با توجه به نتایج هر ۳ دسته نمودار، بیشترین ویژگی‌های انتخاب شده، از ۲ دسته ویژگی CKSAAP و KSCTriad است.



شکل ۳-۱: نمودار ۵۰ ویژگی انتخاب شده توسط انتخاب متوالی رو به جلو



شکل ۳-۲: نمودار ۱۰۰ ویژگی انتخاب شده توسط انتخاب متوالی رو به جلو



شکل ۳-۳: نمودار ۲۰۰ ویژگی انتخاب شده توسط انتخاب متوالی رو به جلو

۳-۴ نتایج بدست آمده بر روی ۳ مجموعه داده SFS50، SFS100 و SFS200

در ادامه نتایج طبقه‌بند جنگل تصادفی بر روی ۳ مجموعه داده SFS50، SFS100 و SFS200 گزارش می‌شود.

۳-۴-۱ نتایج مربوط به طبقه‌بند جنگل تصادفی

نتایج بهترین مدل‌های جنگل تصادفی بر روی داده‌های SFS50، SFS100 و SFS200 در مرحله تنظیم پارامترهای تنظیمی و تست مستقل در جدول زیر گزارش می‌شود.

جدول ۳-۷: نتایج بهترین مدل جنگل تصادفی بر روی داده‌های بدست آمده از الگوریتم SFS

مجموعه داده	صحت (بهترین مدل در تنظیم پارامتر تنظیمی)	معیار roc_auc (بهترین مدل در تنظیم پارامتر تنظیمی)	صحت (تست مستقل)	دقت (تست مستقل)	حساسیت (تست مستقل)	معیار f1 (تست مستقل)
SFS50	۸۲.۰	۸۹.۰	۷۳.۰	۷۳.۰	۷۳.۰	۷۳.۰
SFS100	۸۲.۰	۸۹.۰	۷۴.۰	۷۴.۰	۷۴.۰	۷۴.۰
SFS200	۸۲.۰	۸۹.۰	۷۴.۰	۷۴.۰	۷۴.۰	۷۴.۰

۴ فصل چهارم: بحث و پیشنهادات

با توجه به نتایج طبقه‌بندها در فصل ۳، مشخص است که استفاده از روش‌های محاسباتی برای مسئله پیش‌بینی پتیده‌های ضدسرطانی روش مناسبی است، زیرا نتایج بدست آمده توسط هر ۲ طبقه‌بند جنگل-تصادفی و ماشین بردار پشتیبان با ترکیب ویژگی‌های مختلف و حتی استفاده از یک ویژگی به تنهایی همگی درصد صحت بالاتر از ۷۱ درصد را در تست مستقل گزارش کرده‌اند. همچنین با توجه دقیق‌تر به نتایج طبقه‌بند جنگل تصادفی مشخص می‌شود که استفاده از ویژگی‌های QSOOrder، APseudoAAC، PseudoAAC و CKSAAP بصورت تکی، نتایج بالای ۷۹ تا ۸۱ درصد را در تست مستقل دارند، این در

حالی است که برای طبقه‌بند ماشین بردار پشتیبان فقط برای ویژگی QOrder نتایج تست مستقل بین ۷۹ تا ۸۱ گزارش شد.

در ترکیب ویژگی‌ها نیز طبقه‌بند جنگل تصادفی به نسبت طبقه‌بند ماشین بردار پشتیبان نتایج بالاتری در تست مستقل گزارش می‌کند. در ترکیب ویژگی‌ها، ویژگی‌های QOrder+PseudoAAC و QOrder+APseudoAAC توسط طبقه‌بند جنگل تصادفی با صحت ۸۱ درصد در تست مستقل، بعنوان بیشترین صحت بین ویژگی‌های دیگر گزارش شدند. بالاترین صحت در تست مستقل توسط طبقه‌بند ماشین بردار پشتیبان نیز برای ترکیب ویژگی QOrder+APseudoAAC گزارش شد.

این در حالی است که در مرحله آموزش و تنظیم پارامترهای تنظیمی، صحت بالای ۸۳ درصد در ترکیب ویژگی QOrder+DDE توسط مدل جنگل تصادفی و صحت بالای ۸۲ درصد در ترکیب ویژگی‌های QOrder+CKSAAP، QOrder+APseudoAA و NMBroto+PseudoAAC گزارش شده- است. بالاترین صحت‌های گزارش شده عموماً مربوط به طبقه‌بند جنگل تصادفی است.

از طرفی نتایج صحت بدست آمده توسط جنگل تصادفی بر روی مجموعه داده‌های SFS50، SFS100 و SFS200 نسبت به صحت مجموعه داده‌های حاصل از تحلیل ۱۰ مؤلفه اول بسیار کمتر است طوری که با توجه به جدول ۳-۷ صحت تست مستقل به ۷۳ درصد کاهش یافته‌است. به همین دلیل به نظر می‌رسد که استخراج ویژگی‌های مهم توسط الگوریتم SFS برای پیش‌بینی پتیدهای ضدسرطانی مناسب نیست.

نتایجی که از این پایان‌نامه بدست آمده، نشان می‌دهد که مسأله پیش‌بینی پتیدهای ضدسرطانی مسأله‌ای قابل حل توسط روش‌های محاسباتی با صحت بسیار بالاتر در آینده‌ای نزدیک است. همچنین مشاهده شد

که ویژگی‌های QSOOrder، APseudoAAC، PseudoAAC و CKSAAP ویژگی‌های بهتری نسبت به بقیه ویژگی‌های استخراج شده برای هر دو طبقه‌بند جنگل تصادفی و ماشین بردار پشتیبان هستند.

یکی از چالش‌های حال حاضر برای مسئله پیش‌بینی پپتیدهای ضدسرطانی، کمبود داده‌های آزمایشگاهی پپتیدهای ضدسرطانی است. با دسترسی به داده‌های بیشتر می‌توان به بالا رفتن کارایی مدل‌های یادگیری ماشین و در نتیجه پیش‌بینی‌های با صحت بالاتر کمک کرد. نتیجه پیش‌بینی با صحت بالاتر می‌تواند راه‌های رسیدن به درمان‌های نوین برای سرطان را هموارتر کند. همچنین از بهبودهایی که برای این مطالعه می‌توان در نظر گرفت، می‌توان به مسئله پیش‌بینی پپتیدهای ضدسرطانی برای یک نوع سرطان خاص اشاره کرد؛ که این نیازمند جمع‌آوری داده‌های بیشتر و مخصوصاً جمع‌آوری داده‌های مثبت مربوط به آن نوع سرطان خاص است. مورد دیگری که برای بهبود این مطالعه می‌توان در نظر گرفت، استفاده از روش‌های یادگیری ماشین مدرن (برای مثال استفاده از شبکه‌های عصبی عمیق) و یا استفاده از روش‌های پردازش متن است. که این مورد نیز به دلیل ماهیت شبکه‌های عصبی عمیق، نیازمند داده‌های بسیار بیشتری است.

مراجع

١. Xi, J., et al., *Inferring subgroup-specific driver genes from heterogeneous cancer samples via subspace learning with subgroup indication*. Bioinformatics, 2020. **36**(6): p. 1855-1863.
٢. Yue, Z., et al., *dbCID: a manually curated resource for exploring the driver indels in human cancer*. Briefings in bioinformatics, 2019. **20**(5): p. 1925-1933.
٣. Bray, F., et al., *Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries*. CA: a cancer journal for clinicians, 2018. **68**(6): p. 394-424.
٤. Wijdeven, R.H., et al., *Old drugs, novel ways out: Drug resistance toward cytotoxic chemotherapeutics*. Drug Resistance Updates, 2016. **28**: p. 65-81.
٥. Huang, Y., et al., *Alpha-helical cationic anticancer peptides: a promising candidate for novel anticancer drugs*. Mini reviews in medicinal chemistry, 2015. **15**(1): p. 73-81.
٦. Hoskin, D.W. and A. Ramamoorthy, *Studies on anticancer activities of antimicrobial peptides*. Biochimica et Biophysica Acta (BBA)-Biomembranes, 2008. **1778**(2): p. 357-375.

٧. Giuliani, A., G. Pirri, and S. Nicoletto, *Antimicrobial peptides: an overview of a promising class of therapeutics*. Open Life Sciences, 2007. **2**(1): p. 1-33.
٨. Gaspar, D., A.S. Veiga, and M.A. Castanho, *From antimicrobial to anticancer peptides. A review*. Frontiers in microbiology, 2013. **4**: p. 294.
٩. Ting, C.-H., et al., *The mechanisms by which pardaxin, a natural cationic antimicrobial peptide, targets the endoplasmic reticulum and induces c-FOS*. Biomaterials, 2014. **35**(11): p. 3627-3640.
١٠. Buri, M.V., et al., *Resistance to degradation and cellular distribution are important features for the antitumor activity of gomesin*. PLoS One, 2013. **8**(11): p. e80924.
١١. Vijayakumar, S. and P. Lakshmi, *ACPP: A web server for prediction and design of anti-cancer peptides*. International Journal of Peptide Research and Therapeutics, 2015. **21**(1): p. 99-106.
١٢. Reddy, K., R. Yedery, and C. Aranha, *Antimicrobial peptides: premises and promises*. International journal of antimicrobial agents, 2004. **24**(6): p. 536-547.
١٣. Ejtehadifar, M., et al., *Anti-cancer effects of Staphylococcal Enterotoxin type B on U266 cells co-cultured with Mesenchymal Stem Cells*. Microbial pathogenesis, 2017. **113**: p. 438-444.
١٤. Wang, Z. and G. Wang, *APD: the antimicrobial peptide database*. Nucleic acids research, 2004. **32**(suppl_1): p. D590-D592.
١٥. Schweizer, F., *Cationic amphiphilic peptides with cancer-selective toxicity*. European journal of pharmacology, 2009. **625**(1-3): p. 190-194.
١٦. Wüthrich, K., *NMR with proteins and nucleic acids*. Europhysics News, 1986. **17**(1): p. 11-13.
١٧. Marion, D., M. Zasloff, and A. Bax, *A two-dimensional NMR study of the antimicrobial peptide magainin 2*. FEBS letters, 1988. **227**(1): p. 21-26.

١٨. Ganz, T., et al., *Defensins. Natural peptide antibiotics of human neutrophils*. Journal of Clinical Investigation, 1985. **76**(4): p. 1427.
١٩. Hill, C.P., et al., *Crystal structure of defensin HNP-3, an amphiphilic dimer: mechanisms of membrane permeabilization*. Science, 1991. **251**(5000): p. 1481-1485.
٢٠. Landon, C., et al., *Solution structure of drosomycin, the first inducible antifungal protein from insects*. Protein Science, 1997. **6**(9): p. 1878-1884.
٢١. Tamamura, H., et al., *A comparative study of the solution structures of tachyplesin I and a novel anti-HIV synthetic peptide, T22 ([Tyr 5, 12, Lys 7]-polyphemusin II), determined by nuclear magnetic resonance*. Biochimica et Biophysica Acta (BBA)-Protein Structure and Molecular Enzymology, 1993. **1163**(2): p. 209-216.
٢٢. Xu, T., et al., *Anticandidal activity of major human salivary histatins*. Infection and immunity, 1991. **59**(8): p. 2549-2554.
٢٣. Selsted, M.E., et al., *Indolicidin, a novel bactericidal tridecapeptide amide from neutrophils*. Journal of Biological Chemistry, 1992. **267**(7): p. 4292-4295.
٢٤. Lawyer, C., et al., *Antimicrobial activity of a 13 amino acid tryptophan-rich peptide derived from a putative porcine precursor protein of a novel family of antibacterial peptides*. FEBS letters, 1996. **390**(1): p. 95-98.
٢٥. Gennaro, R., B. Skerlavaj, and D. Romeo, *Purification, composition, and activity of two bactenecins, antibacterial peptides of bovine neutrophils*. Infection and immunity, 1989. **57**(10): p. 3142-3146.
٢٦. AGERBERTH, B., et al., *Amino acid sequence of PR-39*. European Journal of Biochemistry, 1991. **202**(3): p. 849-854.
٢٧. De Vos, W.M., et al., *Properties of nisin Z and distribution of its gene, nisZ, in Lactococcus lactis*. Applied and environmental microbiology, 1993. **59**(1): p. 213-218.

۲۸. Fregeau Gallagher, N.L., et al., *Three-dimensional structure of leucocin A in trifluoroethanol and dodecylphosphocholine micelles: spatial location of residues critical for biological activity in type IIa bacteriocins from lactic acid bacteria*. Biochemistry, 1997. **36**(49): p. 15062-15072.
۲۹. Rokach, L. and O. Maimon, *Decision trees*, in *Data mining and knowledge discovery handbook*. 2005, Springer. p. 165-192.
۳۰. Breiman, L., *Random forests*. Machine learning, 2001. **45**(1): p. 5-32.
۳۱. Osuna, E.E., *Support vector machines: Training and applications*. 1998, Massachusetts Institute of Technology.
۳۲. He, W., et al., *Learning embedding features based on multisense-scaled attention architecture to improve the predictive performance of anticancer peptides*. Bioinformatics, 2021. **37**(24): p. 4684-4693.
۳۳. Chen, Z., et al., *iFeature: a python package and web server for features extraction and selection from protein and peptide sequences*. Bioinformatics, 2018. **34**(14): p. 2499-2502.
۳۴. Hajisharifi, Z., et al., *Predicting anticancer peptides with Chou's pseudo amino acid composition and investigating their mutagenicity via Ames test*. Journal of Theoretical Biology, 2014. **341**: p. 34-40.
۳۵. Chen, W., et al., *iACP: a sequence-based tool for identifying anticancer peptides*. Oncotarget, 2016. **7**(13): p. 16895.
۳۶. Li, F.-M. and X.-Q. Wang, *Identifying anticancer peptides by using improved hybrid compositions*. Scientific reports, 2016. **6**(1): p. 1-6.
۳۷. Manavalan, B., et al., *MLACP: machine-learning-based prediction of anticancer peptides*. Oncotarget, 2017. **8**(44): p. 77121.
۳۸. Schaduangrat, N., et al., *ACPred: a computational tool for the prediction and analysis of anticancer peptides*. Molecules, 2019. **24**(10): p. 1973.

۳۹. Boopathi, V., et al., *mACPpred: a support vector machine-based meta-predictor for identification of anticancer peptides*. International journal of molecular sciences, 2019. **20**(8): p. 1964.
۴۰. van Zoggel, H., et al., *Antitumor and angiostatic activities of the antimicrobial peptide dermaseptin B2*. 2012.
۴۱. Huang, Y., et al., *CD-HIT Suite: a web server for clustering and comparing biological sequences*. Bioinformatics, 2010. **26**(5): p. 680-682.
۴۲. Wei, L., et al., *ACPred-FL: a sequence-based predictor using effective feature representation to improve the prediction of anti-cancer peptides*. Bioinformatics, 2018. **34**(23): p. 4007-4016.
۴۳. Cao, R., et al., *DLFF-ACP: prediction of ACPs based on deep learning and multi-view features fusion*. PeerJ, 2021. **9**: p. e11906.
۴۴. Kawashima, S., et al., *AAindex: amino acid index database, progress report 2008*. Nucleic acids research, 2007. **36**(suppl_1): p. D202-D205.
۴۵. Refaeilzadeh, P., L. Tang, and H. Liu, *Cross-validation*. Encyclopedia of database systems, 2009. **5**: p. 532-538.
۴۶. Raschka, S., *MLxtend: Providing machine learning and data science utilities and extensions to Python's scientific computing stack*. Journal of open source software, 2018. **3**(24): p. 638.
۴۷. James, G., et al., *An introduction to statistical learning*. Vol. 112. 2013: Springer.
۴۸. machinelearningmastery Authors, J.B. *How to Develop a Random Forest Ensemble in Python*. 2020 April 27]; Available from: <https://machinelearningmastery.com/random-forest-ensemble-in-python/>.
۴۹. Towardsdatascience Authors, W.K. *Hyperparameter Tuning the Random Forest in Python*. 2018 Jan 18]; Available from: <https://towardsdatascience.com/hyperparameter-tuning-the-random-forest-in-python-using-scikit-learn-28d2aa77dd74>.
۵۰. ScikitLearn Authors. *sklearn.svm.SVC version 1.0.2*. 2021; Available from: <https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>.

51. DataCamp authors, A.N. *Support Vector Machines with Scikit-learn*. 2019 December 27]; Available from: <https://www.datacamp.com/community/tutorials/svm-classification-scikit-learn-python>.

Abstract

Cancer is one of the major causes of death worldwide. To treat cancer, the use of anticancer peptides (ACPs) has attracted a lot of attention in recent years. ACPs are a unique group of small molecules that can target and kill cancer cells fast and directly. However, identifying ACPs by wet-lab experiments is time-consuming and labor-intensive. Therefore, it is significant to develop computational tools for ACPs prediction. Hence, this study tries to train machine learning models for distinguishing ACPs from non-ACPs. This study utilizes data of ACPs and non-ACPs, uses iFeature python package to extract features from peptide sequences, builds and trains Rrandom Forest and SVM models with scikit-learn python package on extracted features.

Comparing random forest models in this work, random forest with higher than 82% accuracy shows the best performance. Also comparing all models in this work, random forest and SVM models with QSOrder+APseduoAAC features show better performance (accuracy: 81%) on independent test set.

Keywords: Anticancer Peptides Prediction, Machine Learning, Random Forest, SVM, Cancer



Prediction of Anticancer Peptides

Using Machine Learning

Thesis Submitted in Partial Fulfilment of the Requirements for the Degree of Master
of Science (M.Sc.) in Bioinformatics

Department of Biophysics
Faculty of Bioinformatics
Tarbiat Modares University

By:

Mohammadtabar Zeynab

Supervisor:

Dr. Abdolmaleki Parviz

January 2022