

استفاده همزمان از همبستگی خطی پیرسون و ترکیب الگوریتم‌های داده کاوی به منظور بهبود پیش‌بینی نوع تومور در بیماران سرطانی

محسن غلامی^۱، سیدجواد میرعابدینی^{۲*}

۱- دانشجوی کارشناسی ارشد دانشگاه آزاد اسلامی واحد بوشهر. پست الکترونیکی: Mohsen.gholami18@yahoo.com

۲- عضو هیات علمی دانشگاه آزاد اسلامی واحد تهران مرکزی. پست الکترونیکی: j_mirabedini@iauctb.ac.ir

تاریخ دریافت: ۹۷/۱۱/۱۷ تاریخ پذیرش: ۹۸/۵/۱

چکیده

امروزه سرطان سینه از شایع‌ترین بیماری‌های سرطان در بین زنان به‌شمار می‌آید. آمارها از رشد شش درصدی این نوع سرطان در ایران حکایت می‌کند که نشان دهنده جدی بودن خطر آن می‌باشد. این در صورتی است که در صورت پیشگیری و یا تشخیص زود هنگام بیماری می‌توان تا حد زیادی از خطرات آن جلوگیری نمود. با پیشرفت علوم پزشکی، زمینه لازم جهت ایجاد سیستم‌هایی با قابلیت پیشگیری، پیش‌بینی و درمان بیماران با استفاده از فناوری‌های جدید حاصل گردیده است. داده‌کاوی پزشکی سعی در مدل‌سازی و کشف روابط بین عوامل خطر ساز جهت پیش‌بینی وضعیت بیماران آینده با کمک از داده‌های در دست دارد. در این پژوهش سعی گردیده تا با مقایسه الگوریتم‌های مختلف داده‌کاوی و ترکیب این الگوریتم‌ها، روشی جدید، کارا و با دقت بالا و قابلیت پیاده‌سازی بر روی داده‌های محلی ایجاد گردد. در نهایت روش پیشنهادی که به بهبود کارایی الگوریتم بیز ساده با استفاده از الگوریتم آدابوست می‌پردازد، توانایی پیش‌بینی نوع تومور خوش‌خیم یا بدخیم با دقت ۹۶.۶۷ درصد را دارا می‌باشد. داده‌های لازم جهت این فرآیند از سایت UCI جهت تشخیص نوع تومور با ۵۶۹ رکورد و ۳۲ متغیر، استخراج گردیده است.

کلمات کلیدی: ضریب همبستگی پیرسون، الگوریتم‌های دسته‌بندی، بیز ساده، آدابوست

۱- مقدمه :

سرطان به رشد غیر طبیعی سلول‌ها در بدن که به سرعت گسترش می‌یابد گفته می‌شود. طبق تحقیقات صورت گرفته عواملی همچون وراثت، موادشیمیایی و مصرف دخانیات در بروز آن تاثیر گذار هستند [۱].
باتوجه به افزایش سریع روند گسترش بیماری سرطان سینه در زنان و همچنین پایین آمدن سن مبتلایان به این بیماری، لزوم ایجاد سیستم‌هایی جهت تشخیص و پیش‌بینی این بیماری با توجه به اینکه شناسایی زود هنگام عوامل خطر ساز و پیش‌بینی زود هنگام بیماری بسیار در روند بهبود بیماران تاثیر گذار است به شدت احساس می‌گردد چرا که سرطان سینه در مقایسه با سایر سرطان‌ها یکی از عوامل عمده مرگ و میر در میان زنان به‌شمار می‌رود. [۲] از سویی امروزه افزایش هزینه‌ها جهت آزمایشات متعدد و درمان نیز با توجه به گسترش روز افزون بیماری وجود دارد [۳].

از طرف دیگر امکانات پزشکی و تعدد آزمایشات و پزشکان متخصص در همه مناطق جهت تشخیص بیماری وجود ندارد [۴] و داده کاوی پزشکی می‌تواند به کشف روابط پنهان میان متغیرها و ایجاد مدلی توانمند و کارا به سرعت بخشیدن در پیش‌بینی عوامل خطر ساز و پیش‌بینی احتمال عود بیماری به صورت سریع همراه با کاهش هزینه‌ها از طریق کاهش آزمایشات متعدد در تشخیص و پیش‌بینی اولیه بیماری، از طریق تحلیل داده‌های محلی کمک نماید.

طبق تحقیقات صورت گرفته در کشورهای اروپایی سن ابتلا به سرطان سینه در زنان بین سنین ۵۳ الی ۵۷ سال است [۵]. این در صورتی است که میانگین سن ابتلا به این بیماری در ایران تقریباً ۴۵ سال می‌باشد. از این رو ضرورت ایجاد سیستم‌هایی جهت تشخیص عوامل خطر ساز و پیش‌بینی زود هنگام بیماری در مراحل اولیه بیماری امری ضروری به نظر می‌رسد [۶].

ایجاد سیستم‌های ساده و موثر جهت کاهش تعداد متغیرها و افزایش دقت و کارایی الگوریتم‌ها یکی از شاخص‌های حیاتی برای تشخیص و پیش‌بینی زود هنگام و کاهش هزینه‌های آزمایشات متعدد، یکی از اهداف اصلی داده‌کاوی پزشکی محسوب می‌گردد [۷]. حتی این سیستم با استفاده از داده‌های مرتبط، می‌تواند در تشخیص نوع تومور مثل خوش‌خیم بودن و یا بدخیم بودن آن به پزشکان متخصص کمک کند.

در این مقاله سعی شده ابتدا الگوریتم‌های مختلف از نظر دقت و کارایی با هم مقایسه گردیده و سپس با روش‌های ترکیبی کارایی این الگوریتم‌ها بهبود یافته و در نهایت روشی کارا و ساده ایجاد می‌گردد. روش پیشنهادی علاوه بر کاهش تعداد متغیرها، باعث افزایش دقت پیش‌بینی در مقایسه با الگوریتم اولیه گردیده و عملیاتی می‌باشد. از این روش نیز می‌توان به تشخیص نوع تومور از نظر خوش‌خیمی و بدخیمی به متخصصین یاری رساند.

۲- ضرورت انجام تحقیق:

امروزه سرطان یکی از علل اصلی مرگ و میر در جهان به‌شمار می‌آید و تشخیص و پیش‌بینی زود هنگام، در مرحله رشد بیماری و تشخیص نوع آن یکی از چالش‌های جدی در علوم پزشکی است. که می‌توان با پیش‌بینی و تشخیص زود هنگام و به موقع، آن را تا حد زیادی کنترل و از بروز آن جلوگیری کرد [۸].

سرطان سینه به عنوان شایع‌ترین سرطان در بین زنان می‌باشد و روند رو به رشد آن در سال‌های اخیر بسیار زیاد و نگران‌کننده بوده است. پیش‌بینی احتمال بروز این بیماری و یا تشخیص نوع توده و یا احتمال عود مجدد، می‌تواند تا حد زیادی در بهبود و یا کنترل بیماری تاثیرگذار باشد [۹].

کاهش سن مبتلایان به سرطان سینه در کشور که به میانگین سنی ۴۵ سال رسیده است در صورتی که در کشورهای غربی این نوع از سرطان بین زنان ۵۳ الی ۵۷ سال شایع است.

در تحقیقات صورت گرفته در خصوص علل مرگ و میر و روند تغییرات آن در سال‌های ۱۳۵۸ تا ۱۳۸۰ در ایران که توسط دکتر پروین یاری و همکاران صورت پذیرفت نتایج حاصله در بازه ۲۳ ساله نشان می‌داد مرگ به علت بیماری‌های واگیردار در کشور کاهش یافته ولی متاسفانه بیماری‌های غیر واگیر مثل بیماری‌های قلبی و سرطان سیر صعودی داشته‌اند.

از طرفی ایجاد مدل‌های پیش‌بینی و تشخیص می‌توانند به تسهیل روند مدیریت و شناسایی بیماری کمک نمایند [۱۰].

۳- مفاهیم اولیه

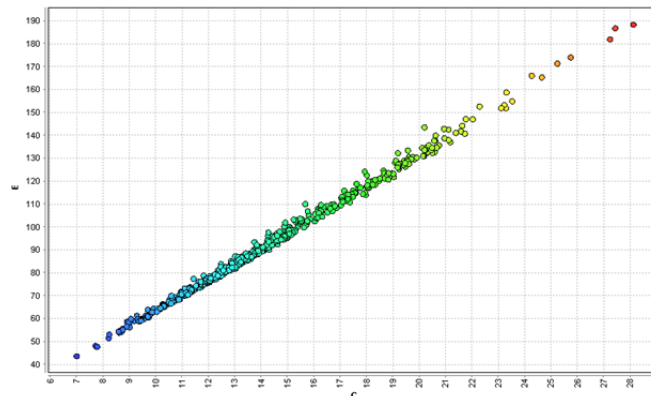
۱-۳ ضریب همبستگی پیرسون

با استفاده از این تکنیک آماری می‌توان میزان همبستگی خطی بین دو متغیر تصادفی را شناسایی نمود. این ضریب مقداری بین ۱- و ۱ دارد. همبستگی کامل زمانی رخ می‌دهد که ضریب ما ۱ باشد و زمانی که ضریب ۰ باشد یعنی همبستگی بین متغیرها وجود ندارد. ضریب ۱- نیز نشان دهنده نسبت منفی کامل بین دو مقدار است [۱۱]. در واقع ما با استفاده از این تکنیک میزان ارتباط خطی بین متغیرها را شناسایی می‌کنیم که از فرمول زیر قابل محاسبه است:

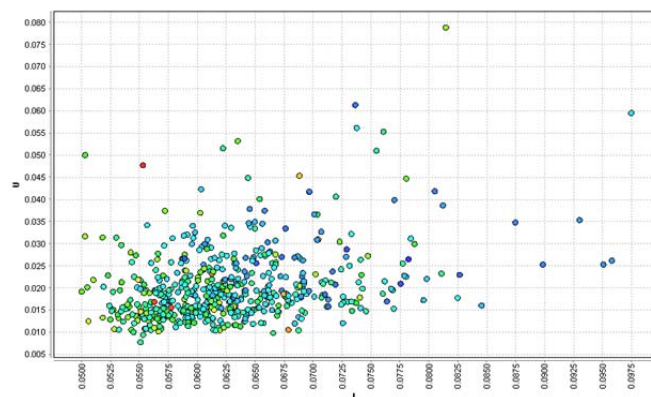
$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

فرمول شماره (۱) ضریب همبستگی جامعه

فرمول شماره (۲) ضریب همبستگی نمونه آماری پیرسون



شکل (۱) پراکندگی قرارگیری داده‌ها بر روی محور مختصات که نشان می‌دهد همبستگی خطی زیادی بین دو متغیر E و C وجود دارد.



شکل (۲) نشان می‌دهد با توجه به قرارگیری داده‌ها به صورت دایره‌ای همبستگی خطی کمی بین دو متغیر L و U وجود دارد (نزدیک به صفر)

۲-۳ واریانس و کواریانس:

واریانس نحوه پراکندگی داده‌ها در اطراف میانگین همان داده‌ها را مشخص می‌کند همچنین کواریانس شاخصی است برای تغییرات یک متغیر با متغیر دیگر است [۱۲].

۳-۳ الگوریتم بیز ساده

در یادگیری ماشین معمولاً در فضای فرضیه H بدنال بهترین فرضیه‌ای هستیم که درمورد داده‌های آموزشی D صدق کند. یک راه تعیین بهترین فرضیه، این است که بدنال محتمل‌ترین فرضیه‌ای باشیم که با داشتن داده‌های آموزشی D و احتمال قبلی در مورد فرضیه‌های مختلف می‌توان انتظار داشت تئوری بیز چنین راه‌حلی را ارائه می‌دهد. این روش راه حل مستقیمی است که نیازی به جستجو ندارد [۱۳].

سنگ‌بنای یادگیری بیزی را تئوری بیز تشکیل می‌دهد. این تئوری امکان محاسبه احتمال ثانویه را بر مبنای احتمالات اولیه می‌دهد [۱۴].

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

فرمول شماره (۳) قضیه بیز

همان‌طور که مشاهده می‌شود با افزایش $P(D)$ مقدار $P(h|D)$ کاهش می‌یابد. زیرا هر چه احتمال مشاهده D مستقل از h بیشتر باشد به این معنا خواهد بود که D شواهد کمتری در حمایت از h در بر دارد. [۱۵] [۱۶]

قضیه بیز:

$P(h)$ = احتمال اولیه‌ای که فرضیه h قبل از مشاهده مثال آموزشی D داشته است [۱۶].

$P(D)$ = احتمال اولیه‌ای که داده آموزشی D مشاهده خواهد شد [۱۶].

$P(D|h)$ = احتمال مشاهده داده آموزشی D به فرض آنکه فرضیه h صادق باشد [۱۶].

۴-۳ الگوریتم آدابوست^۱

آدابوست مخفف بوستینگ تطبیقی بوده و یک الگوریتم یادگیری ماشین است که توسط یاو فروند و رابرت شاپیر ابداع شد [۱۷]. در واقع آدابوست یک متا الگوریتم است که بمظور ارتقاء عملکرد، و رفع مشکل رده‌های نامتوزان [۱۸] همراه دیگر الگوریتم‌های یادگیری استفاده می‌شود. در این الگوریتم، طبقه‌بند هر مرحله جدید به نفع نمونه‌های غلط طبقه‌بندی شده در مراحل قبل تنظیم می‌گردد [۱۹]. آدابوست نسبت به داده‌های نویزی و پرت حساس است؛ ولی نسبت به مشکل بیش‌برازش از بیشتر الگوریتم‌های یادگیری برتری

^۱- Adaboost

دارد. طبقه‌بند پایه که در اینجا استفاده می‌شود فقط کافیست از طبقه‌بند تصادفی (۵۰٪) بهتر باشد و به این ترتیب بهبود عملکرد الگوریتم با تکرارهای بیشتر بهبود می‌یابد [۲۰].

۳-۵ جنگل تصادفی ۱

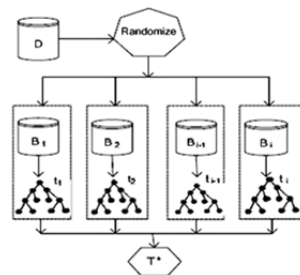
جنگل تصادفی مجموعه‌ای از درختان تصمیم است که بصورت موازی با هم در ارتباط اند و بر پایه تئوری بگینگ تشکیل شده است [۲۰] [۲۱]. این روش دسته‌بندی یادگیری با ناظر است که در آن درخت‌های تصمیم صفات خاصه خود را بطور تصادفی از داده‌های کاملاً مستقل انتخاب و سپس دسته‌بندی می‌کند.

جنگل تصادفی را می‌توان با دو روش Forest-RI و Forest-RC ایجاد کرد [۲۲]. در روش Forest-RI که بر پایه روش bagging است جنگل تصادفی بصورت تصادفی و پی در پی صفات خاصه را انتخاب و در نهایت دسته‌بندی می‌کند [۲۳]. برای ساخت و رشد درختان در جنگل تصادفی می‌توان از روش درخت تصمیم CART استفاده کرد.

در روش Forest-Rc ابتدا با استفاده از ترکیب تعدادی از صفات خاصه و تعیین ضرائب در محدوده (1, -1) صفات خاصه را تصادفی انتخاب و در نهایت بهترین صفت را گسترش می‌دهد. از معایب این روش، محدودیت تعداد صفات خاصه است. در صورتی که تعداد صفات خاصه زیاد شود باعث کاهش همبستگی میان دسته‌بندها می‌شود [۲۴].

جنگل تصادفی امکان هرس کردن درختان در صورت بزرگ شدن را ندارد و به همین دلیل نرخ خطا برای جنگل تعریف می‌شود. صحت و درستی یک جنگل تصادفی به قدرت و توانایی هر یک از دسته‌بندها و وابستگی میان آن‌ها بستگی دارد.

از مزایای جنگل تصادفی می‌توان به کارایی و دقت بالا نسبت به روش‌های Bagging و Boosting اشاره کرد. این الگوریتم بر روی داده‌های بزرگ با ویژگی‌های زیاد بدون حذف متغیر نتایج خوبی را ارائه می‌دهد [۲۵]. همچنین در برخورد با داده‌های پرت و نامتوازن بهتر از سایر الگوریتم‌ها جهت دسته‌بندی عمل می‌کند [۲۶].



شکل ۳ نحوه عملکرد الگوریتم جنگل تصادفی [۲۷].

روش کار این الگوریتم به صورتی است که ابتدا k نمونه با استفاده از الگوریتم Bagging به صورت تصادفی انتخاب می‌شود. در مرحله بعد، از k نمونه‌ها برای k درخت و k دسته‌بند استفاده می‌شود. در نهایت از بین دسته‌بندی کننده‌ها برای دسته‌بندی بهینه رای گیری صورت می‌گیرد [۲۸].

۳-۶ درخت تصمیم

الگوریتم درخت تصمیم، یک الگوریتم بازگشتی دارای ساختار درختی است که از ریشه به سمت برگ حرکت می‌کند. این الگوریتم یکی از روش‌های طبقه‌بندی داده بر پایه یادگیری با نظارت است. گره‌های برگ حاوی اطلاعات کلاس‌ها (متغیر وابسته) و گره‌های غیر برگ صفات خاصه (متغیر مستقل) می‌باشد. روش کار الگوریتم بصورت سلسله مراتبی است که بر اساس داده‌های آموزشی با انتخاب یکی از صفات خاصه در هر مرحله شروع به کار می‌کند. در ادامه با تقسیم‌بندی هر یک از صفات ادامه می‌دهد تا زمانیکه تمام داده‌ها به اطلاعات دارای برچسب واحد کلاس شوند. این روش ساده و به راحتی قابل درک و تفسیر بوده و امکان هرس کردن تصمیم‌هایی که قابل تعمیم نیستند را دارد. الگوریتم‌های درخت تصمیم به صورت چند مرحله‌ای عمل می‌کنند یعنی تصمیمات پیچیده به تصمیمات ساده تری تقسیم و ترکیب این تصمیمات ساده رسیدن به تصمیمات مورد انتظار است [۲۹].

معیارهای مختلفی جهت انتخاب صفات خاصه در این الگوریتم وجود دارد که می‌توان به موارد زیر اشاره نمود:

Information Gain: این معیار مشخص کننده این است که یک ویژگی تا چه مقدار می‌تواند مثال‌های آموزشی را بر اساس دسته بندی تفکیک کند [۲۴]. این معیار در ID3 برای انتخاب صفات خاصه استفاده می‌شود.

$$IG(A) = \text{Entropy}(D) - \text{Entropy}(A|D)$$

فرمول ۴) محاسبه بهره اطلاعاتی

$$\text{Entropy}(D) = -\sum_{i=1}^c p_i \times \log_2(p_i)$$

فرمول ۵) محاسبه آنتروپی

این فرمول‌ها برای محاسبه مفهوم آنتروپی که در واقع میزان خلوص (بی نظمی یا عدم خالص بودن) مجموعه‌ای از مثال‌ها را مشخص می‌کند [۲۰] بکار برده می‌شود. در توضیح فرمول شماره ۱ و ۲ بهتر است بگوییم:

■ **C** معرف تعداد برچسب‌های کلاس موجود در داده‌ها.

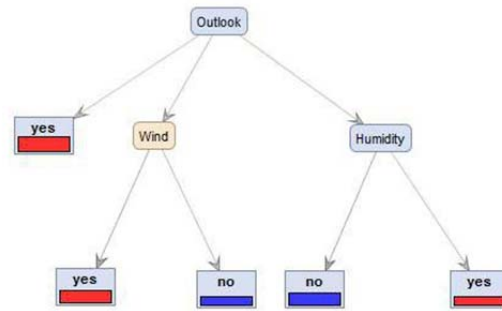
■ **Pi** معرف احتمال اینکه نمونه‌ای از داده‌ها متعلق به کلاس i ام باشد.

■ **V** تعداد اعضای دامنه صفت خاصه A .

■ **Dj** قسمتی از داده اولیه که مقدار صفت خاصه آن‌ها V_j باشد.

Gain Ratio: این معیار در واقع برای نرمال‌سازی **Information Gain** استفاده می‌شود. این روش صفات خاصه دارای دامنه زیاد را به صفات خاصه با دامنه کم ترجیح می‌دهد که دارای دقت بالا و همچنین کاهش پیچیدگی نسبت به عملکرد **Information Gain** است.

Gini Index: این معیار تمام صفات خاصه را بصورت دودویی دسته بندی می‌کند سپس مقدار **Gini** را محاسبه کرده و کمترین مقدار صفت خاصه را انتخاب می‌کند. این معیار برای انتخاب گره در الگوریتم **CART** مورد استفاده قرار می‌گیرد [۳۰].



شکل ۴) نمونه ای از رشد درخت از ریشه به برگ

در توضیح شکل ۱، **Outlook** گره ریشه و **wind** و **Humidity** گره‌ها و مقادیر **yes** و **no** برگ‌ها هستند. به طور کلی از مزایای درخت تصمیم می‌توان به امکان استفاده برای انواع داده‌های پیوسته و گسسته، استخراج قوانین قابل فهم و آسان، آماده‌سازی آسان داده‌ها، سازگار کردن داده‌های فاقد مقدار، شناسایی تفاوت میان زیرگروه‌ها، دسته‌بندی ویژگی‌ها با تأثیرات زیاد و محاسبات کم برای دسته‌بندی داده‌ها و تفسیر آسان درخت اشاره کرد. معایب این الگوریتم، مصرف بالای حافظه، سخت و زمانبر بودن پیاده‌سازی برای مجموعه داده‌های بزرگ، هزینه بالا برای هرس کردن، زیاد بودن تعداد گره‌های پایانی در صورت همپوشانی، انباشته شدن خطای لایه‌ها بر روی یکدیگر در صورت بزرگ شدن درخت است.

۳-۷ K نزدیکترین همسایه ۱:

این الگوریتم برای اولین بار در سال ۱۹۵۰ معرفی شد. اما به دلیل سرعت پایین محاسبات کامپیوترها و نیاز به حافظه بالا در آن زمان، تا سال ۱۹۶۰ مورد توجه قرار نگرفت [۳۰]. یادگیری خود را بر اساس میزان تشابه و فاصله انجام می‌دهد. از این الگوریتم در تشخیص الگو و پیش‌بینی بصورت گسترده استفاده می‌شود.

در این روش دسته‌بندی داده‌ها بصورت بردار یا نقاطی در فضای چند بعدی ویژگی‌ها به گونه‌ای که مفهومی به نام فاصله بوجود می‌آید صورت می‌پذیرد. بعد از آن هریک از داده‌ها نرمال‌سازی شده چراکه به داده‌های نویز حساس و در صورت وجود بر روی صحت پاسخگویی اثر می‌گذارد. در این الگوریتم می‌توان از فاصله اقلیدوسی یا فاصله منتهن برای سنجش فاصله میان صفات استفاده کرد. اگر بطور کلی هر یک از مجموعه داده‌ها را در دو کلاس مثبت و منفی برچسب بزنیم، این الگوریتم با در نظر گرفتن k تعداد دلخواه همسایه را برای تعیین و انتخاب کلاس‌ها برای دسته‌بندی انتخاب می‌نماید که بیشتر توسط کاربر انتخاب می‌شود. پیچیدگی ذخیره‌سازی و محاسباتی این الگوریتم $O(n)$ است [۱۱].

این الگوریتم دارای دقت خوبی است. سرعت انجام محاسبه آن پایین است. یکی از مشکلات این روش این است که قابلیت تعمیم ندارد (برخلاف روش رگرسیون) و با اضافه کردن داده جدید می‌بایست دوباره مدل سازی صورت گیرد. این الگوریتم را **Lazy** یا تنبل نیز می‌نامند چرا که تا داده آزمایشی وارد نشود مدل سازی انجام نمی‌گردد.

برای محاسبه فاصله می‌توان از معیارهای متفاوتی استفاده کرد که یکی از این معیارها فاصله اقلیدوسی است. فاصله دو نقطه p و q اندازه پاره‌خطی است که آنها را به هم متصل می‌کند که می‌توان به صورت زیر تعریف کرد.

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2} = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

فرمول (۶) محاسبه فاصله اقلیدسی

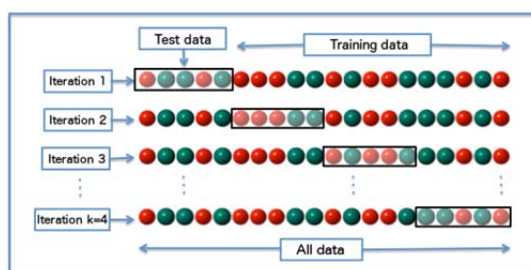
۳-۸ الگوریتم RIPPER

از این الگوریتم برای کشف قوانین به صورت مستقیم با استفاده از معیار بهره اطلاعاتی استفاده می شود. این الگوریتم شامل دو مرحله رشد و هرس کردن است. هرس کردن تا زمانیکه میزان خطا بیشتر از ۵۰ درصد باشد ادامه می یابد. به عبارت ساده تر قوانینی که تاثیری بر روی دقت مدل ندارد ضعیف شناخته شده و حذف می شود.

۳-۹ روش ارزیابی مدل با استفاده از K-Fold

در این پژوهش از روش K-Fold جهت ارزیابی مدل استفاده شده است. در این روش نوع اعتبارسنجی داده ها به K زیرمجموعه افراز می شوند. از این K زیرمجموعه، هر بار یکی برای اعتبارسنجی و K-1 تای دیگر برای آموزش بکار می روند. این روال K بار تکرار می شود و همه داده ها دقیقاً یکبار برای آموزش و یکبار برای اعتبارسنجی بکار می روند [۳۱]. در نهایت میانگین نتیجه این K بار اعتبارسنجی به عنوان یک تخمین نهایی برگزیده می شود. البته می توان از روش های دیگر برای ترکیب نتایج استفاده کرد. بطور معمول از ۱۰-Fold استفاده می شود [۳۲].

در روش K-Fold طبقه ای سعی می شود نسبت داده های هر کلاس در هر زیرمجموعه و در مجموعه اصلی یکسان باشد.



شکل (۵) نحوه انتخاب مجموعه داده تست جهت ارزیابی در روش K-Fold [۳۳].

۴- داده ها

در این مقاله از مجموعه داده استاندارد سرطان سینه ویسکانسین از مخزن یادگیری UCI استفاده می گردد که حاوی ۵۶۹ پرونده و دارای ۳۲ متغیر است که جهت تشخیص نوع تومور استفاده می شود و متغیر هدف آن دارای دو نوع خوش خیم و بدخیم می باشد. لازم به توضیح است مقادیر ویژگی ها از تصاویر عکس برداری شده استخراج شده اند [۳۴].

جدول ۱ معرفی ویژگی‌های مجموعه داده [۳۵]

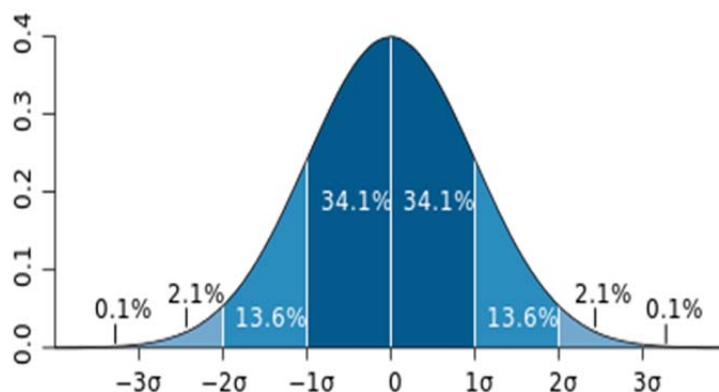
ردیف	نام متغیر	توضیحات
۱	ID	شماره پرونده
۲	diagnosis	متغیر هدف که شامل دو مقدار خوش خیم و بد خیم است
۳	radius	شعاع
۴	texture	بافت
۵	perimeter	محیط
۶	area	مساحت
۷	smoothness	همواری
۸	compactness	میزان غلظت
۹	concavity	میزان فرورفتگی
۱۰	concave_points	نقاط توخالی (تعداد مقادیر توخالی حد فاصل)
۱۱	symmetry	تقارن
۱۲	Fractal_dimension	ابعاد فراکتال
۱۳	radius_SE	شعاع SE
۱۴	texture_SE	بافت SE
۱۵	perimeter_SE	محیط SE
۱۶	area_SE	مساحت SE
۱۷	smoothness_SE	همواری SE
۱۸	compactness_SE	فشرده‌گی SE
۱۹	concavity_SE	فرورفتگی SE
۲۰	concave_points_SE	نقاط مقعر SE
۲۱	symmetry_SE	تقارن SE
۲۲	fractal_dimension_SE	ابعاد فراکتال SE
۲۳	radius_Worst	بدترین شعاع
۲۴	texture_Worst	بدترین بافت
۲۵	perimeter_Worst	بدترین محیط
۲۶	area_Worst	بدترین مساحت
۲۷	smoothness_Worst	بدترین همواری
۲۸	compactness_Worst	بدترین فشرده‌گی
۲۹	concavity_Worst	بدترین تورفتگی
۳۰	concave_points_Worst	نقاط تو رفتگی ارزشمند
۳۱	symmetry_Worst	بدترین تقارن
۳۲	fractal_dimension_Worst	بدترین بعد فراکتال

۴-۱ پیش پردازش داده ها:

۴-۱-۱ بررسی مقادیر خارج از محدوده و گم شده

انحراف معیار اندازه گیری چگونگی پراکندگی یک ویژگی است که با کمک آن می توان داده های ناسازگار (با توجه به میزان فاصله) را مشخص نمود. انحراف به معنی دوری از مقدار نرمال است. انحراف معیار، عددی برای نشان دادن میزان گسترش اعداد است. علامت آن حرف یونانی سیگما σ است. با استفاده از انحراف معیار، ما یک راه استاندارد برای یافتن مقدار نرمال، مقدار بیش از نرمال و مقدار کمتر از نرمال در دست داریم.

یک قانون سرانگشتی خوب این است که معمولاً داده های بین $\bar{x} - 2\sigma$ و $\bar{x} + 2\sigma$ داده های با ارزش تری محسوب شده و داده های خارج از این فاصله در نظر گرفته نمی شود که اصطلاحاً به داده های کم ارزش داده های خارج از محدوده^۱ نیز می گویند [۳۲]. می توان با استفاده از نمودار ۳-۱ این قضیه را بهتر درک کرد که به آن اصطلاحاً نمودار توزیع نرمال داده های تصادفی یا نمودار محدوده های قابل نرمال گویند.



نمودار ۱: نمودار توزیع طبیعی [۱۱]

در توضیح نمودار ۱ باید بگوییم که قسمت آبی تیره در فاصله یک برابر انحراف معیار از میانگین توزیع قرار دارد و قسمت آبی روشن و آبی تیره به طور توأم، در فاصله دو برابر انحراف معیار از میانگین توزیع قرار دارند. در توزیع طبیعی، اولی برابر با ۶۸٪ سطح زیر نمودار و دومی برابر با ۹۵٪ سطح زیر نمودار است.

معمولاً با افزایش تعداد داده ها توزیع آن ها به منحنی توزیع نرمال میل پیدا می کند [۱۱]. در توزیع نرمال، ۶۸.۲٪ داده ها در فاصله کمتر از یک انحراف معیار نسبت به میانگین قرار دارند. این مقدار برای فاصله های دو و سه انحراف معیار، به ترتیب ۹۵.۴٪ و ۹۹.۷٪ است. به بیان دیگر، احتمال آن که اختلاف یک داده با میانگین، بیش از سه انحراف معیار باشد، تنها ۰.۳٪ (تقریباً معادل ۱ در ۳۰۰) است [۱۱] [۳۶].

برای مثال، صفت age که میانگین آن ۳۶.۷۳۱ است در حالی که انحراف معیار آن ۱۰.۶۴۷ است. حال مطابق با قانون سر انگشتی معرفی شده می بایست فاصله دو برابر انحراف معیار از میانگین را محاسبه کنیم که به صورت زیر عمل می کنیم:

^۱- outlier

$$(36.731 + (2 * 10.647)) = 58.025$$

$$(36.731 - (2 * 10.647)) = 15.437$$

فرمول شماره (۷): محاسبه فاصله دو برابر انحراف معیار از میانگین

صفت age دارای یک بازه از ۱۷ تا ۵۷ است که نشان می‌دهد همه مشاهدات ما در فاصله کمتر از دو برابر انحراف معیار از میانگین که ۵۸.۰۲۵ و ۱۵.۴۳۷ است قرار دارد و در محدوده می‌باشد.

در بررسی‌های صورت گرفته مشخص گردید مقادیر مجموعه داده ما فاقد مقادیر خارج از محدوده است و گم شده است.

۵- روش کار مدل پیشنهادی

در این روش ابتدا متغیرهای دارای ضریب همبستگی بیشتر با استفاده از ضریب همبستگی پیرسون^۱ شناسایی شده و در ادامه با انتخاب متغیرهای قوی‌تر (بالتر از ۰.۴+ یا پایین‌تر از ۰.۴-) مدل‌سازی صورت می‌گیرد. با مقایسه نتایج حاصل از مدل‌سازی این روش انتخاب متغیر کارایی بهتری نسبت به سایر روش‌های انتخاب صفات از خود نشان می‌دهد و بیست متغیر جهت مدل‌سازی انتخاب می‌شود.

در ادامه از الگوریتم بیز ساده جهت مدل‌سازی استفاده شده و با استفاده از الگوریتم آدابوست تقویت شده و توانایی لازم جهت پیش-بینی نوع تومور را از خود نشان می‌دهد. در این پژوهش آدابوست وظیفه دارد تا هفت دسته بند متفاوت برای بیز ساده ایجاد کند که بیشتر روی تاپل‌های سخت‌تر تمرکز دارد منظور از تاپل‌های سخت‌تر مقادیری است که در دست بندهای قبلی غلط پیش‌بینی شده اند. آدابوست با انتخاب زیر مجموعه‌هایی از مجموعه داده اصلی به مدل‌سازی پرداخته و در نهایت هفت دسته بند ایجاد شده به وسیله رای‌گیری و در نهایت نیز برای پیش‌بینی از همین هفت دسته بند استفاده می‌شود.

با توجه به هدف این تحقیق، که بهبود کارایی و دقت پیش‌بینی نوع تومور بیماران سرطانی است، ابتدا چندین الگوریتم دسته‌بندی مقایسه شده و پس از آن با استفاده از ضریب همبستگی پیرسون، متغیرهایی که دارای میانگین همبستگی و ارتباط بیشتری نسبت به سایر متغیرها هستند انتخاب و عمل آموزش و ایجاد مدل جدید بر روی این متغیرها صورت می‌گیرد. با حذف متغیرهای ضعیف که دارای همبستگی صفر یا کمی هستند علاوه بر افزایش کارایی پارامترهایی همچون Accuracy, Recall, Precision, Kappa تقویت شده و میزان بازه احتمالی خطا در هر بار آموزش نیز کاهش پیدا کرده است و در ادامه با استفاده از الگوریتم آدابوست، فرآیند آموزش بهبود پیدا می‌کند. وجود بازه احتمالی خطا به دلیل این است که در این پژوهش از روش Cross Validation جهت ارزیابی مدل استفاده گردیده است. دقت مدل پیشنهادی پس از ارزیابی ۹۶.۶۷ درصد می‌باشد.

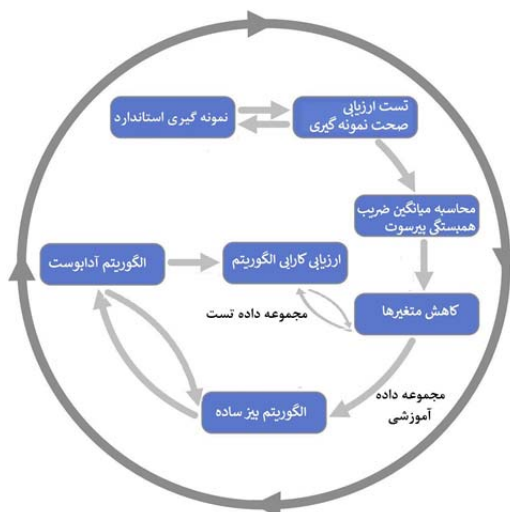
۵-۱ مراحل اجرایی

مراحل اجرایی روش پیشنهادی به شرح زیر است:

۱. نمونه‌گیری استاندارد
۲. تست ارزیابی صحت نمونه‌گیری

^۱ Pearson correlation

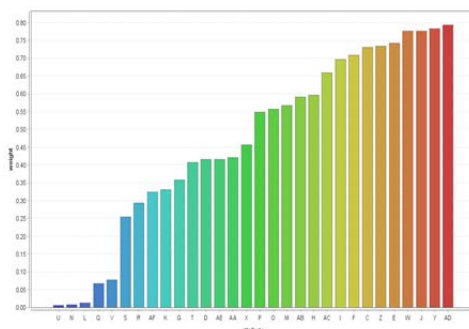
۳. محاسبه میانگین ضریب همبستگی با روش پیرسون و انتخاب متغیرهای قوی تر
۴. کاهش ابعاد و حذف متغیرهای با همبستگی کمتر
۵. فرآیند آموزش با الگوریتم بیز ساده
۶. ترکیب الگوریتم بیز ساده^۱ با الگوریتم آدابوست
۷. تست مدل ترکیبی ضریب همبستگی با پیرسون و بیز ساده و آدابوست



شکل ۶) فلوجارت روش پیشنهادی

۶- کاهش تعداد متغیرها

جهت بهبود پارامترهای مختلف با استفاده از میانگین ضریب همبستگی پیرسون برای متغیرها کاهش ابعاد صورت می گیرد و بیست متغیر برتر انتخاب می شود و عملیات آموزش و تست صورت می گیرد.



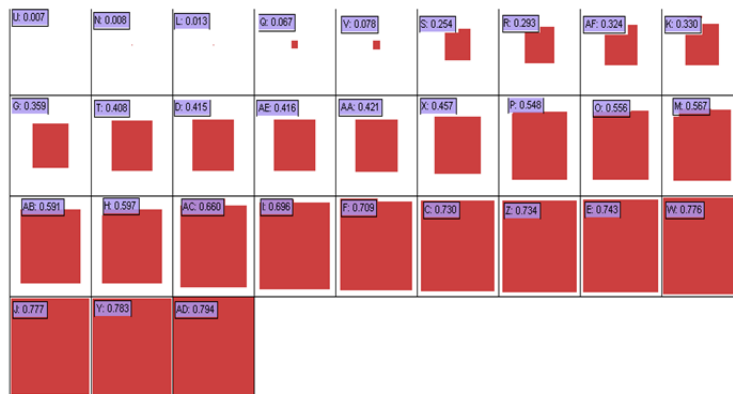
نمودار ۲: تست و مرتب سازی میزان همبستگی متغیرها با روش پیرسون

¹ Naive Bayes

استفاده همزمان از همبستگی خطی پیرسون و ترکیب الگوریتم‌های داده کاوی به منظور

جدول ۲: مقدار وزن تعیین شده به همراه نام متغیر

وزن	نام متغیر	ردیف
۰.۷۹۳۵۶۶	AD	۱
۰.۷۸۲۹۱۴	Y	۲
۰.۷۷۶۶۱۴	J	۳
۰.۷۷۶۴۵۴	W	۴
۰.۷۴۲۶۳۶	E	۵
۰.۷۳۳۸۲۵	Z	۶
۰.۷۳۰۰۲۹	C	۷
۰.۷۰۸۹۸۴	F	۸
۰.۶۹۶۳۶	I	۹
۰.۶۵۹۶۱	AC	۱۰
۰.۵۹۶۵۳۴	H	۱۱
۰.۵۹۰۹۹۸	AB	۱۲
۰.۵۶۷۱۳۴	M	۱۳
۰.۵۵۶۱۴۱	O	۱۴
۰.۵۴۸۲۳۶	P	۱۵
۰.۴۵۶۹۰۳	X	۱۶
۰.۴۲۱۴۶۵	AA	۱۷
۰.۴۱۶۲۹۴	AE	۱۸
۰.۴۱۵۱۱۵	D	۱۹
۰.۴۰۸۰۴۲	T	۲۰
۰.۳۵۸۵۶	G	۲۱
۰.۳۳۰۴۹۹	K	۲۲
۰.۳۲۳۸۷۲	AF	۲۳
۰.۲۹۲۹۹۹	R	۲۴
۰.۲۵۳۷۳	S	۲۵
۰.۰۷۷۹۷۲	V	۲۶
۰.۰۶۷۰۱۶	Q	۲۷
۰.۰۱۲۸۳۸	L	۲۹
۰.۰۰۸۳۰۳	N	۲۹
۰.۰۰۶۵۲۲	U	۳۰



نمودار ۳: نمودار Hinton ضریب همبستگی متغیرها

در توضیح نمودار ۳ هر مربع نشان دهنده یک متغیر است و هر چه مربع قرمز رنگ بزرگ تر باشد، نشان دهنده ضریب همبستگی بیشتر آن است.

۷- شبه کد روش پیشنهادی

Calculate Pearson correlation coefficient

$$(1) \Gamma = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$$

Selecting variables with more correlation coefficient

Algorithm: AdaBoost. A boosting algorithm-Create an ensemble of classifiers. Each one gives a weighted vote.

Input:

- D, a set of d class- labled training tuples;
- K, the number of round (one classifier is generated per round);
- a classification learning scheme.

Output: A composite model.

Method:

- (1) initialize the weight of each tuple in D to 1/d;
- (2) for i=1 to k do // for each round:
- (3) sample D with replacement according to the tuple weights to obtain D_i;
- (4) Use training set D_i to derive a model, M_i;
- (5) Compute error(M_i), the error rate of M_i (Eq. 8.34)
- (6) If error(M_i) > 0.5 then
- (7) go back to step 3 and try again;
- (8) Endif
- (9) For each tuple in D_i that was correctly classified do
- (10) Multiply the weight of the tuple by error(M_i)/(1-error(m_i)); //update weights
- (11) Normalize the weight of each tuple;

(12) Endfor

To use the ensemble to classify tupe, X:

- (1) Initialize weight of each class to 0;
- (2) For $i=1$ to k do //for each classifier:
- (3) $W_i = \log \frac{1 - \text{error}(M_i)}{\text{error}(M_i)}$; // weight of the classifiers vote
- (4) $C = M_i(X)$; // get class prediction for X from M_i
- (5) Add w_i to weight for class c
- (6) Endfor
- (7) Return the class with the largest weight;

۱-۷ تشریح کدهای به کار رفته

استفاده از فرمول پیرسون ضریب همبستگی میان متغیرها شناسایی شده و متغیرهای دارای ضریب بیشتر از $+0.4$ و کمتر از -0.4 انتخاب می‌شود.

$$\Gamma = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$$

فرمول شماره (۸): محاسبه ضریب همبستگی خطی پیرسون

در این فرمول X و Y دو متغیر و n تعداد تاپل‌ها است. به عنوان مثال در مجموعه داده D با دو متغیر X و Y در صورتی که پرونده ۶ بیمار وجود داشته باشد مقدار n عدد ۶ خواهد بود. همچنین $\sum xy$ ضرب مقادیر هر ردیف متغیر X در Y خواهد بود که در نهایت با هم جمع می‌شود. هدف از این کار انتخاب متغیرهای قوی‌تر و حذف متغیرهای ضعیف‌تر است، که در نهایت بیست متغیر جهت مدل‌سازی مرحله بعد انتخاب می‌شود. در ادامه مجموعه داده که دارای کاهش تعداد متغیرها است با روش آدابوست تقویت می‌شود. آدابوست الگوریتمی از شاخه بوسستینگ است.

در روش **boosting** به هر یک از تاپل‌های آموزشی وزنی تخصیص داده می‌شود [20]. پس از ساخته دسته‌بند m_i وزن‌ها تغییر خواهند کرد تا دسته‌بند m_{i+1} که پس از m_i تولید می‌شود، توجه بیشتری را بر روی تاپل‌هایی که به درستی توسط m_i دسته‌بندی نشده‌اند، داشته باشند.

دسته‌بندی نهایی m^* رای نهایی را با ترکیب رای‌های هر یک از دسته‌بندهای پایه محاسبه می‌کند. جایی که وزن رای هر یک از این دسته‌بندها تابعی از صحت آن است.

از میان روش‌های **boosting** الگوریتم **Adaboost** یکی از الگوریتم‌های رایج محسوب می‌شود که در این پژوهش نیز از این روش، برای افزایش یادگیری مدل استفاده می‌شود. مجموعه داده‌های D از تعداد d تاپل که هر یک به شکل (x_i, y_i) هستند تشکیل شده است و در آن بر چسب کلاس تاپل x_i با y_i نشان داده می‌شود.

در ابتدا الگوریتم Adaboost مقدار وزن یکسان $1/d$ را برای هر تاپل آموزشی در نظر می‌گیرد که در خط شماره ۱ شبه کد نشان داده شده است. d تعداد تاپل‌ها است. در این روش تلفیقی برای تولید k دسته‌بند به k مرحله نیاز است که ما تعداد k را عدد ۷ در نظر گرفته‌ایم. در این پژوهش از الگوریتم بیز ساده به عنوان دسته‌بند پایه استفاده می‌شود.

در مرحله i ام با کمک نمونه‌گیری با جایگزینی مجموعه آموزشی D_i شکل می‌گیرد.

در این روش ممکن است تاپل‌های یکسانی بیش از یک بار انتخاب شود. شانس انتخاب m تاپل بر اساس وزن تعیین می‌شود.

با استفاده از مجموعه داده‌های آموزشی D_i مدل m_i ساخته می‌شود و با کمک همین مجموعه، خطای مدل نیز محاسبه می‌گردد. پس از آن وزن‌های تاپل‌های آموزشی بر اساس اینکه چگونه دسته‌بندی شده‌اند، تنظیم می‌شود.

اگر تاپلی به درستی دسته‌بندی نشده باشند، وزن آن افزایش می‌یابد و وزن تاپل‌هایی که به درستی دسته‌بندی شده‌اند کاهش داده می‌شود.

در واقع وزن یک تاپل، دشواری دسته‌بندی آن تاپل را منعکس می‌کند، به صورتی که وزن بالاتر نشان می‌دهد اغلب این تاپل‌ها به درستی دسته‌بندی نشده است.

وزن‌های جدید تاپل‌ها برای تولید مجموعه آموزشی دسته‌بند بعدی در مرحله‌ی بعد استفاده می‌شود. ایده اصلی کار از جایی سرچشمه می‌گیرد که ما هنگام ساخت دسته‌بند مایلیم که مدل بیشتر بر روی دسته‌بندی تاپل‌هایی تمرکز کند که در مرحله قبلی به درستی دسته‌بندی نشده‌اند.

ممکن است برخی از دست‌بندها در دسته‌بندی تاپل‌های دشوار بهتر از بعضی دیگر عمل کنند، بدین طریق یک سری از دسته‌بندهایی ساخته شده‌اند که همدیگر را تکمیل می‌کنند.

مجموع وزن‌های تاپل‌هایی در D_i که مدل M_i آنها را به درستی دسته‌بندی نکرده است، به عنوان نسبت خطای مدل در نظر گرفته می‌شود، بنابراین داریم:

$$\text{error}(M_i) = \sum_{j=1}^d w_j \times \text{err}(X_j)$$

فرمول شماره (۹): محاسبه خطای دسته بند

که در آن تابع $\text{err}(X_j)$ هرگاه تاپل X_j به درستی دسته‌بندی نشده باشد، برابر با ۱ و در غیر این صورت برابر با صفر خواهد بود. چنانچه خطای حاصل از دسته‌بند M_i از مقدار ۰/۵ تجاوز کند، کارایی مدل ضعیف تشخیص داده می‌شود و از مدل صرف نظر خواهد شد، در عوض با تولید یک مجموعه آموزشی جدید D_i سعی در ساختن مدل جدید M_i خواهیم داشت. نسبت خطای M_i بر روی چگونگی بهنگام‌سازی وزن‌های تاپل‌های آموزشی تاثیر می‌گذارد. چنانچه تاپلی در مرحله i ام به درستی دسته‌بندی شود، وزن آن در مقدار $(1 - \text{error}(M_i)) / \text{error}(M_i)$ ضرب می‌شود (خط ۱۰ شبه کد). پس از بهنگام‌سازی وزن کلیه تاپل‌هایی که به درستی دسته‌بندی شده‌اند، وزن تمام تاپل‌ها (حتی آن‌هایی که به درستی دسته‌بندی نشده‌اند) به نحوی نرمال‌سازی می‌شود که مجموع وزن‌ها همانند قبل باقی بماند. برای نرمال‌سازی یک وزن کافی است آن را در مجموع وزن‌های قدیم ضرب و بر مجموع وزن‌های جدید تقسیم کنیم.

در نتیجه همانطور که قبل از این نیز توضیح داده شد، با این کار وزن تاپل‌هایی که به درستی دسته‌بندی نشده‌اند افزایش و وزن تاپل‌هایی که به درستی دسته‌بندی شده‌اند کاهش می‌یابد. پس از پایان الگوریتم، از روش تلفیقی دسته‌بندها برای پیشگویی برچسب کلاس تاپلی مانند X استفاده می‌شود. بر خلاف روش بگینگ که در آن تمام دسته‌بندها دارای رای یکسانی هستند، در این روش بر اساس عملکرد هر دسته‌بند وزنی برای رای آن در نظر گرفته می‌شود. در واقع رای‌های آن‌ها دارای وزن‌های یکسانی نیست. هر چه نسبت خطای مدل کمتر باشد، در نتیجه صحت آن بیشتر خواهد بود و بنابراین وزن بالاتری را برای رای آن لحاظ می‌کنیم. وزن رای دسته بند $M_i =$ برابر است با:

$$\log \frac{1 - \text{error}(M_i)}{\text{error}(M_i)}$$

فرمول شماره (۱۰): محاسبه وزن رای دسته بند

مجموع وزن‌های دسته‌بندی‌هایی که کلاس C را برای X پیشگویی می‌کنند محاسبه می‌کنیم. این کار برای تمام برچسب کلاس‌ها انجام می‌شود و بیشترین مقدار، تعیین کننده کلاس تاپل X خواهد بود.

۷-۲ محاسبه وزن رای دسته بندها:

همانطور که در تشریح مدل ذکر گردید در الگوریتم آدابوست در نهایت با استفاده از رای گیری دسته‌بندها پیش‌بینی انجام می‌گردد. البته این رای‌گیری به این صورت است که با توجه به وزن هر دسته‌بند با توجه به صحت آن مشخص می‌شود. این باعث می‌شود که تاثیر رای دسته‌بندها با هم یکسان نباشد. (هر چه وزن دسته‌بند بیشتر باشد، تاثیر رای بیشتری دارد). برای محاسبه وزن دسته بند از فرمول زیر استفاده می‌کنیم.

$$\log \frac{1 - \text{error}(M_i)}{\text{error}(M_i)}$$

فرمول شماره (۱۱): محاسبه نسبت خطای دسته بند

$\text{Error}(M_i)$ همانطور که قبلاً نیز تشریح داده شد نسبت خطای هر دسته‌بند است.

در مدل پیشنهادی از ۷ دسته‌بند استفاده شده که وزن هر دسته‌بند به صورت زیر می‌باشد. در تشریح این مقادیر بهتر است بگوییم Model 1 نسبت خطای کمتر و در نتیجه صحت بیشتری نسبت به دسته‌بندهای دیگر داشته است در نتیجه رای آن تاثیر بیشتری نسبت به سایر دسته‌بندها دارد.

Adaboost

Model 1 [w= 2.820] (Naive Bayes)

Model 2 [w= 1.333] (Naive Bayes)

Model 3 [w= 1.495] (Naive Bayes)

Model 4 [w= 1.231] (Naive Bayes)

Model 5 [w= 0.273] (Naive Bayes)

Model 6 [w= 0.587] (Naive Bayes)

Model 7 [w= 0.975] (Naive Bayes)

۳-۷ روش ارزیابی مدل و تاثیر افزایش حجم داده‌ها بر روی مدل

در این پژوهش از روش Cross validation (k-fold) جهت ارزیابی مدل استفاده گردیده است. این تکنیک بطور ویژه در کاربردهای پیش‌بینی مورد استفاده قرار می‌گیرد تا مشخص شود مدل موردنظر تا چه اندازه در عمل مفید خواهد بود. این روش مناسب برای ارزیابی مجموعه داده‌هایی با تعداد کم و برای دوری از Overfitting پیشنهاد می‌شود. در این مقاله از مقدار پیش‌فرض و متعارف تعداد تکرار مدل و مجموعه داده که ده بار است، استفاده شده است. نتایج حاصل از بررسی‌ها با استفاده از این روش نشان می‌دهد مدل پیشنهادی قابلیت تعمیم در دنیای واقعی را داشته و توانایی لازم در صورت افزایش تعداد بیماران و حجم داده‌ها را با سرعت مناسب دارد.

لازم به توضیح است الگوریتم‌ها با لپ‌تاپ Lenovo G50 با CPU: Core i3 ایتمل اجرا و مدل‌سازی شده‌اند.

همچنین برای محاسبه کارایی روش پیشنهادی از نظر سرعت و مواجه با تعداد بالای بیماران ابتدا با ۵۶۹ پرونده سپس با ۱۱۳۸ پرونده و بعد از آن با ۱۷۰۷ پرونده و در نهایت با ۲۲۷۶ پرونده مدل سازی گردید که روش پیشنهادی در بیشترین تعداد بیماران قادر به پاسخگویی در کمتر از ۴ ثانیه بود که نسبت به سایر مدل‌ها بسیار سریع‌تر است.

۸- نتایج مدل سازی روش پیشنهادی

با استفاده از روش ترکیبی کاهش ابعاد با استفاده از اندازه‌گیری میزان همبستگی توسط تکنیک پیرسون و حذف متغیرهایی که دارای همبستگی کمتری هستند و استفاده از الگوریتم بیز ساده و ترکیب آن با الگوریتم آدابوست این الگوریتم بهینه‌سازی شده و کارایی و دقت آن به نحو قابل توجهی افزایش پیدا می‌کند.

جدول ۳: ماتریس نتایج مدل سازی روش پیشنهادی

	true M	true B	class precision
pred. M	۲۰۳	۱۰	%۹۵.۳۱
pred. B	۹	۳۴۷	%۹۷.۴۷
class recall	%۹۵.۷۵	%۹۷.۲۰	
Kappa Kohen	۰.۹۲۸		

در توضیح جدول ۳ قطر اصلی نشان دهنده تعداد پیش‌بینی صحیح مدل و قطر فرعی نشان دهنده تعداد خطای مدل است. به عنوان مثال مدل پیشنهادی ۳۴۷ مورد را (Benign) B و ۲۰۳ مورد را (Malignant) M به درستی پیش‌بینی کرده است. از این ماتریس برای محاسبه پارامترهای مختلف مدل استفاده می‌گردد.

جدول ۴: مدل میزان کارایی روش پیشنهادی

RD	Algorithm	Accuracy	Recall	Precision	Kappa Kohen
1	Proposed method	%۹۶.۶۷ +/-۲.۷۷	%۹۶.۴۶ +/-۳.۶۱	%۹۶.۵۳ +/-۲.۴۱	۰.۹۲۸

در توضیح جدول ۴ معیار Kappa هر چه مقدارش نزدیک به عدد یک باشد عملکرد بهتری دارد. ارزیابی معیارهای مختلف نشان دهنده کارایی بالای روش پیشنهادی است. مقادیر بدست آمده در این جدول حاصل از ماتریس ارزیابی مدل جدول شماره ۳ است به صورتی که پس از محاسبه مقدار دقت مدل پیشنهادی 96.67 درصد تخمین زده شده است که برای محاسبه آن از فرمول زیر استفاده می‌کنیم.

$$\text{Accuracy} = (203 + 347) / 569 = 96.67$$

فرمول شماره (۱۲): محاسبه مقدار پارامتر Accuracy

مقادیر فوق از جدول شماره ۳ که حاصل از ارزیابی مدل است استخراج گردیده است. قطر اصلی این ماتریس که شامل اعداد ۲۰۳ و ۳۴۷ است تعداد بیمارانی است که سیستم به درستی پیش‌بینی کرده و به آن برچسب درست داده است. تعداد کل مجموعه داده نیز ۵۶۹ مورد است. از این رو دقت مدل ۹۶.۶۷ درصد تخمین زده شده است که بالاتر از سایر روش‌ها است.

۸-۱ مقایسه روش پیشنهادی با سایر روش‌ها:

جدول ۵: مقایسه روش پیشنهادی با سایر روش‌ها

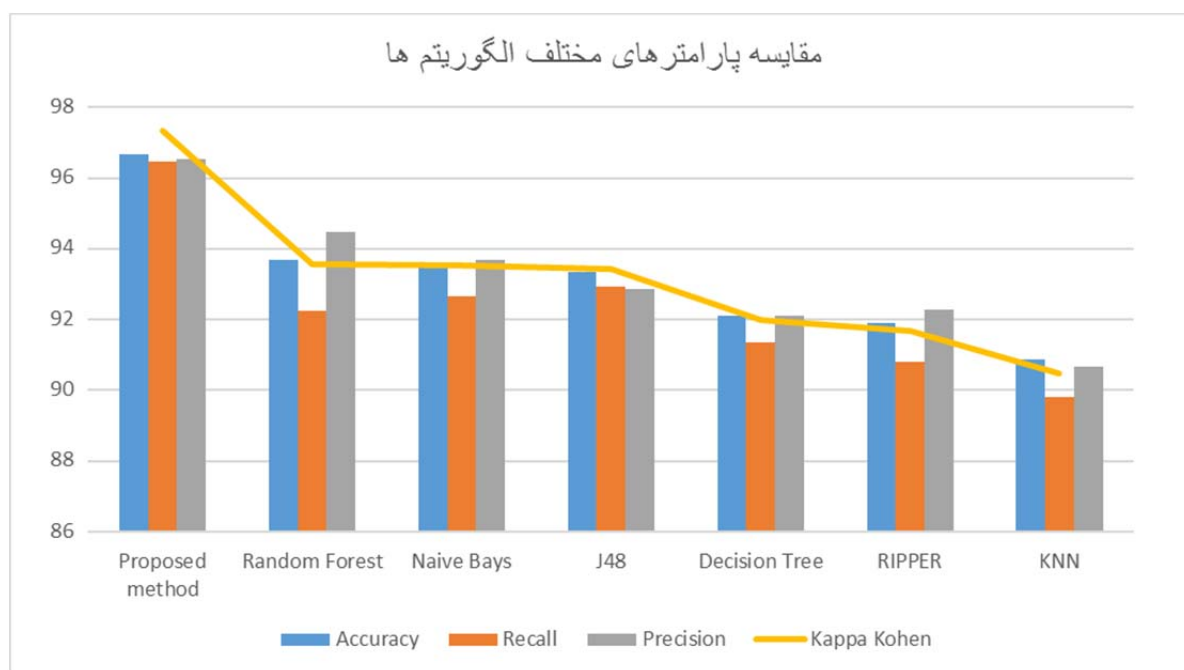
RD	Algorithm	Accuracy	Recall	Precision	Kappa Kohen
۱	Proposed method	%۹۶.۶۷ +/-۲.۷۷	%۹۶.۴۶ +/-۳.۶۱	%۹۶.۵۳ +/-۲.۴۱	۰.۹۲۸
۲	Random Forest	%۹۳.۶۸ +/-۵.۲۷	%۹۲.۲۴ +/-۶.۷۹	%۹۴.۴۸ +/-۴.۳۱	۰.۸۵۹
۳	Naive Bays	%۹۳.۵۱ +/-۵.۰۲	%۹۲.۶۷ +/-۶.۲۴	%۹۳.۶۸ +/-۴.۵۰	۰.۸۵۸
۴	J48	%۹۳.۳۳ +/-۴.۰۶	%۹۲.۹۳ +/-۴.۹۰	%۹۲.۸۷ +/-۴.۰۵	۰.۸۵۶
۵	Decision Tree	%۹۲.۱۰ +/-۳.۹۴	%۹۱.۳۷ +/-۵.۲۳	%۹۲.۱۱ +/-۳.۳۹	۰.۸۳۰
۶	RIPPER	%۹۱.۹۲ +/-۴.۳۱	%۹۰.۷۹ +/-۵.۰۵	%۹۲.۲۹ +/-۳.۹۹	۰.۸۲۴
۷	KNN	%۹۰.۸۶ +/-۳.۲۱	%۸۹.۸۲ +/-۴.۱۶	%۹۰.۶۶ +/-۲.۹۹	۰.۸۰۲

نتایج جدول شماره ۵ نشان می‌دهد، مدل پیشنهادی ترکیبی از پیرسون، بیز و آداپوست در همه پارامترهای ارزیابی شده از سایر الگوریتم‌ها بهتر عمل می‌کند. توانایی این روش ترکیبی وقتی بهتر مشخص می‌گردد که در کنار پارامترها با توجه به استفاده از روش CrossValidation میانگین خطای مدل در ده بار ارزیابی به ۲.۷۷ درصد کاهش پیدا کرده است در صورتی که مدل اولیه بیز ساده تلورانس ۵.۰۲ درصد داشت. مدل پیشنهادی نیز دارای دقت ۹۶.۶۷ درصد است. پس از مدل پیشنهادی الگوریتم درخت تصادفی با

دقت ۹۳.۶۸ درصد قرار می‌گیرد این الگوریتم دارای خطای میانگین ۵.۲۷ درصدی در ده بار مدل سازی است که نسبت به روش ترکیبی بسیار بیشتر است.

۸-۲ نقاط قوت الگوریتم پیشنهادی نسبت به سایر الگوریتم‌ها

دلیل تقویت مدل پیشنهادی نسبت به سایر الگوریتم‌ها این است که تکنیک آدابوست بر خلاف سایر روش‌ها بر روی مقادیری که پیش‌بینی آن‌ها سخت‌تر است تمرکز دارد. به این صورت که پس از مدل‌سازی اولیه سیستم با ترکیب مجموعه داده‌های مختلف در کنار یکدیگر سعی دارد تا قوانینی را کشف و هویدا سازد که بتواند داده‌های سخت را به درستی برچسب‌دهی کند از این رو با وزن دار کردن آن‌ها میزان تمرکز خود را بر روی داده‌ها افزایش یا کاهش می‌دهد. طبیعی است که داده‌های سخت‌تر بیشتر در مدل‌سازی شرکت کرده تا در نهایت سیستم توانایی پیش‌بینی وضعیت آن‌ها را نیز داشته باشد در صورتی که این تکنیک در سایر الگوریتم‌های بکار گرفته شده کاربرد ندارد.



نمودار ۴: مقایسه کارایی سایر الگوریتم‌ها با روش پیشنهادی

۹- نتیجه‌گیری

با توجه به روند روبه رشد سرطان سینه در بین زنان و پایین آمدن سن مبتلایان به این بیماری در ایران، لزوم تشخیص نوع تومور به صورت خوش خیم و یا بدخیم، با دقت بالا به شدت احساس می‌شود. در این پژوهش با استفاده از ضریب همبستگی پیرسون، متغیرهایی که دارای میانگین ضریب همبستگی بیشتری هستند، انتخاب شده و در ادامه متغیرهای با همبستگی کمتر از مجموعه داده حذف گردید و کاهش ابعاد صورت گرفت. مقایسه عملکرد الگوریتم‌ها قبل و بعد از کاهش ابعاد، نشان می‌دهد دو الگوریتم بیزساده

و شبکه‌های عصبی، به‌نحو قابل توجهی در پارامترهای مختلف بهبود یافته‌اند. از این‌رو در ادامه برای بهینه‌سازی بیشتر الگوریتم بیز ساده با الگوریتم آدابوست نیز ترکیب شد. حاصل ترکیب استفاده از ضریب همبستگی پیرسون و بیز ساده با آدابوست، ایجاد مدلی ترکیبی است که توانایی پیش‌بینی و شناسایی نوع تومور سرطانی با دقت ۹۶.۶۷ درصد را دارا است. از مزایای روش پیشنهادی سهولت پیاده‌سازی و سرعت و دقت بیشتر می‌توان اشاره نمود. استفاده از همبستگی پیرسون و روش بکار رفته در مقایسه با روش‌های پیشین راحت‌تر با آموزش ساده‌تر و با دقت بالاتر می‌باشد.

منابع

- [1] Thangaraju, P., Mehala, R., (2015); "Novel Classification based approaches over Cancer Diseases", International Journal of Advanced Research in Computer and Communication Engineering, Vol 4, Issue 4, 294-297
- [2] Karim Khani Zand Hamid, (2015); "A COMPARITIVE SURVEY ON DATA MINING TECHNIQUES FOR BREAST CANCER DIAGNOSIS AND PREDICTION", Indian Journal of Fundamental and Applied Life Sciences, Vol ۵, ۴۳۳۹-۴۳۳۰
- [3] Majidi Zolbanin Hamed, Delen Dorsan, Hassan Zadeh Amir, (201۵); "Predicting overall survivability in comorbidity of cancers: A data mining approach", Elsevier Decision Support Systems 74,150-161
- [4] Venkatalakshmi B, Shivsankar M.V, (2014); "Heart Disease Diagnosis Using Predictive Data mining", International Journal of Innovative Research in Science, Engineering and Technology, Volume 3, Issue 3,1873-1877
- [5] Vikas Chaurasia, Saurabh Pal, (2014); "A Novel Approach for Breast Cancer Detection using Data Mining Techniques", International Journal of Innovative Research in Computer and Communication Engineering, Vol 2, Issue 1,2456-2465
- [6] غلامی، م. برومندیان، ع. (۱۳۹۵)، ارائه روشی با استفاده از ترکیب بهره اطلاعاتی، K نزدیکترین همسایه و شبکه‌های عصبی جهت پیش‌بینی وضعیت جنین، یازدهمین کنفرانس علوم و تکنولوژی، مشهد، ایران، بهمن ۹۵.
- [7] میرعابدینی سیدجواد، غلامی، م. (۱۳۹۵)؛ «ارائه روشی جهت افزایش دقت تشخیص بیماری‌های قلبی با ترکیب الگوریتم‌های داده‌کاوی (MAWB)»، سومین همایش ملی مهندسی رایانه و مدیریت فناوری اطلاعات، سوم، تهران، دانشگاه شهید بهشتی
- [8] P.Ramachandran, Dr.N.Girija, Dr.T.Bhuvanewari, (2013); "Cancer Spread Pattern – an Analysis using Classification and Prediction Techniques", International Journal of Advanced Research in Computer and Communication Engineering, Vol 2, Issue 6,2363-2367
- [9] Kharya Shweta, (2012); "USING DATA MINING TECHNIQUES FOR DIAGNOSIS AND PROGNOSIS OF CANCER DISEASE", International Journal of Computer Science, Engineering and Information Technology (IJCSEIT), Vol 2, Issue 2,55-66
- [10] Safdari Reza, Ghazisaedi Marjan, et al (2013); "A Model for Predicting Myocardial Infarction Using Data Mining Techniques", Iranian Journal of Medical Informatics, Vol 2, Issue 4,1-6
- [11] غلامی، م. (۱۳۹۶)، داده‌کاوی برای همه، تهران، انتشارات ناقوس، چاپ اول.
- [12] North, M., (2012), Data mining for the masses. Amazon, First Edition
- [13] Purusothaman, G., Krishnakumari, P., (2015), A Survey of Data Mining Techniques on Risk Prediction: Heart Disease, Indian Journal of Science and Technology, Vol 8(12), 2-5.
- [14] غلامی، م. نجفی، ن. (۱۳۹۵)، بررسی و مقایسه اثربخشی الگوریتم‌های داده کاوی جهت پیش‌بینی بیماری پارکینسون، اولین کنفرانس بین‌المللی چشم اندازه‌های نو در مهندسی برق و کامپیوتر، تهران، ایران، دانشگاه علم و صنعت، بهمن ۹۵.
- [15] Durgalakshmi, B., Vijayakumar, V., (2015), Prognosis and Modelling of Breast Cancer and its Growth Novel Naive Bayes, Procedia Computer Science, Vol 50, 551, 553.

- [16] چوبینه، پ. غلامی، م. (۱۳۹۵)، مقایسه شش الگوریتم برتر حوزه داده کاوی، یازدهمین کنفرانس علوم و تکنولوژی، مشهد، ایران، بهمن ۹۵.
- [17] Guru Rao, C.V., Sreenivasa Rao, M., (2016), Cluster Analysis of Medical Research Data using R, Global Journal of Computer Science and Technology: C Software & Data Engineering, Vol 16, 17–22.
- [18] Karim Khani Zand, H., (2015), A COMPARITIVE SURVEY ON DATA MINING TECHNIQUES FOR BREAST CANCER DIAGNOSIS AND PREDICTION, Indian Journal of Fundamental and Applied Life Sciences, Vol 5, 4330-4339.
- [19] Boughorbel, S., Al-Ali, R., Elkum, N., (2016), Model Comparison for Breast Cancer Prognosis Based on Clinical Data, PLOS ONE, 15-1.
- [20] Han, J., Jian, P., (2011), Data mining: concepts and techniques, Elsevier.
- [21] 5. L. Breiman, J. H. Friedman. (1984), "Classification and regression trees," Monterey.
- [22] حقیقی. مهری، "داده کاوی و یادگیری ماشین: مروری بر دسته بندی کننده ها"، کنفرانس بین المللی یافته های نوین، تهران، ایران، شهریور ۹۴.
- [23] مرتضی پور. رضا، مطلق زاده، مهسان، استفاده از شبکه بیزین در طبقه بندی، چهارمین کنفرانس مهندسی برق و الکترونیک، دانشگاه آزاد گناباد، شهریور، ۱۳۹۱.
- [24] نظری. احمد، "ارائه الگوریتم ترکیبی بهینه بردار ماشین پشتیبان و جنگل تصادفی در تشخیص به موقع بیماری های قلبی"، کنفرانس بین المللی پژوهش در علوم و تکنولوژی، مالزی، اذر ۹۴.
- [25] Y. Liu and P. Zhao and et al. (2015), "A Boosting Algorithm for Item Recommendation with Implicit Feedback," Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence (IJCAI 2015)
- [26] Krishnaiah, V., Narsimha, G., Subhash Chandra, N., (2016), Heart Disease Prediction System using Data Mining Techniques and Intelligent Fuzzy Approach: A Review, International Journal of Computer Applications, Vol 136, No 2, 43–51.
- [27] سعادت. محمد، زمانی بروجنی. فرساد، "مروری بر روش های بهبود کارایی تکنیک جنگل تصادفی"، اولین همایش ملی فناوری اطلاعات، ارتباطات و محاسبات نرم، دانشگاه آزاد خوراسگان، اصفهان، ایران. اردیبهشت ۱۳۹۵.
- [28] Han, Y. Liu, X. Sun. (2013), "A scalable random forest algorithm based on MapReduce," Software Engineering and Service Science (ICSESS), 2013 4th IEEE International Conference on, pp.849-852.
- [29] Parijaee Moghaddam. A and Mousavi. S , Learning Decision Tree Using Neural Network for Stability and Flexibilit, Iranian Journal of Medical Informatics, vol. 1, no. 3, pp. 39-44, 2013.
- [30] T. M. Khoshgoftar, M. Golawala, and J. Van Hulse. (2007), "An Empirical Study of Learning from Imbalanced Data Using Random Forest.," presented at the 19th IEEE Conference on Tools with Artificial Intelligence.
- [31] Zheng, G. (2017), Logistic Regression, Model Selection, and Cross Validation, personal.umich.edu, March. 25.
- [32] Soundarya, M., Balakrishnan, R., (2014), Survey on Classification Techniques in Data mining, International Journal of Advanced Research in Computer and Communication Engineering, Vol. 3, 7550-7552.
- [33] Zheng, G. Logistic Regression, Model Selection, and Cross Validation, personal.umich.edu, March. 25, 2017.
- [34] Street, W., Wolberg, W., Mangasarian, O., (1992), Nuclear Feature Extraction For Breast Tumor Diagnosis. International Symposium on Electronic Imaging Science and Technology, VOL, 1905, 861-870.
- [35] Lichman, M. (2013). UCI Machine Learning Repository [http://archive.ics.uci.edu/ml]. Irvine, CA: University of California, School of Information and Computer Science.
- [36] Christian, S., Winston, W., Zappa, Ch., (2011), Data Analysis and Decision Making, Forth edition, 14-16.

use Pearson's Linear Correlation and the combination of Data Mining Algorithms simultaneously to improve prognosis of a kind of tumor in cancer patients

mohsen gholami¹, Seyed Javad Mirabedini^{2*}

Islamic Azad University Bushehr, Iran - mohsen.gholami18@outlook.com
Department of Computer, Central Tehran Branch, Islamic Azad University, Tehran, Iran j_mirabedini@iauctb.ac.ir

Abstract:

Nowadays, breast cancer is the most common cancer disease among women. Statistics shows a six percent increase in Iran which indicates it as a serious danger. However, its danger can be prevented increasingly by early diagnosis or prediction. By medical science progress, the way for developing of a system with the capability of prevention, prognosis and cure by using the new technologies is paved. Medical data mining tries to design a model and find relationships among risky factors to predict the condition of future patients with the aid of current data. We try to compare different data mining algorithms and combination of these algorithms to develop a new, efficient method with high accuracy and capability to perform on local data. Finally, proposed method which improves efficiency of Naive Bayes with Adaboost algorithm can predict the kind of benign or malign tumor with the 96/67% accuracies. Required data for this procedure is extracted from UCI site to diagnose the kind of tumor with 569 records and 32 variables.

Keywords: Pearson's correlation coefficient, classification algorithms, Naive Bayes, Adaboost