

معرفی و ارزیابی میزان دقیقیت تکنیک‌های داده کاوی برای کمک به پیش‌بینی بیماری‌های قلبی

محمد متوجهی^۱ و سمية فرهنگ‌آدیب^۲

^۱دانشجوی کارشناسی ارشد، دانشگاه قم، Arjmand.2200@yahoo.com

^۲کارشناسی ارشد، دانشگاه قم، Somayehadib@yahoo.com

چکیده - امروزه با افزایش حجم داده‌ها در سیستم‌های اطلاعاتی، نیاز به ابزاری است تا با کمک آن بتوان داده‌های ذخیره شده را مورد بردازش قرارداده و اطلاعات حاصل از این بردازش را در اختیار کاربران قرار داد. تکنیک‌های داده کاوی یکی از مهمترین و بهترین روش‌ها است که به وسیله آن الگوهای مفید از داده‌ها خارج می‌شود و پس از تجزیه و تحلیل مورد استفاده قرار می‌گیرد. یکی از کاربردهای اصلی داده کاوی در علم پژوهشی مخصوصاً در شاخه بیماری‌های قلبی است که از مهمترین علل های مرگ و میر در جهان می‌باشد. در این تحقیق از چندین روش کلاسی بندی متداول نظیر روش ماشین بردار پشتیبان، درخت تصمیم، بیزین و چندین روش دیگر استفاده شده است و در آخر یک روش ترکیبی ابداعی نیز معرفی شده است. این روش‌ها بر روی مجموعه داده‌های بیماری‌های قلبی که از ترکیب داده‌های Clementine ۴ مرکز پژوهشی مختلف جمع آوری شده است به منظور ارزیابی میزان دقیق و مقایسه روش‌ها با همیگر پیاده شده است. برای محاسبه دقیق این روش‌ها از ابزار Clementine استفاده شده است. نتایج حاصل از این روش‌ها بر روی هر ۴ مجموعه داده نشانگر این است که بهترین الگوریتم بستگی به ساختار داده‌ها، داشت اولیه و مسئله مورد نظر دارد.

dataacademy.ir

کلید واژه- تکنیک‌های داده کاوی، بیماری‌های قلبی، داده کاوی، ابزار Clementine. روش ترکیبی

کنونی تا سال ۲۰۲۵ میلادی بیشتر از ۳۵ تا ۶۰ درصد موارد مرگ و میر در جهان از بیماری‌های قلبی عروقی ناشی شود.^[۱]

۱- مقدمه

استفاده از روش‌های آماری در بررسی داده‌ها مدتی است که رونق گرفته است، به طوری که اغلب تحلیل‌های پژوهشی به وسیله روش‌های آماری نظیر SPSS انجام می‌شود. روش‌های آماری معمولاً برای اثبات فرضیه مورد نظر استفاده می‌شود، بدین ترتیب که ابتدا نظریه‌ای طرح شده و سپس توسط روش‌های آماری، درستی آن مورد ارزیابی قرار می‌گیرد. همکاری متخصصان در زمینه کامپیوتر و پژوهشی هستند نه کشف و کاربردی تحلیل این داده‌ها و به دست آوردن الگوهای مفید و کاربردی ارائه می‌دهد که همان داده کاوی است.^[۲] در داده کاوی برخلاف علم آمار به دنبال پیشگویی هستند نه کشف یا اثبات. بدین معنا که با استفاده از روش‌های داده کاوی به دنبال تأیید آنچه از قبل وجود دارد نیستند بلکه به دنبال مشخص کردن الگوهای از قبل شناخته نشده هستند.^[۳] همچنین در داده کاوی بیماری‌های قلبی به دنبال شناسایی افرادی که دارای بیماری قلبی

امروزه در داشت پژوهشی جمع آوری داده‌های فراوان در مورد بیماری‌های مختلف از اهمیت فراوانی برخوردار است. مراکز پژوهشی با مقاصد گوناگونی به جمع آوری این داده‌ها می‌پردازن. تحقیق روی این داده‌ها و به دست آوردن نتایج و الگوهای مفید در رابطه با بیماری‌ها، یکی از اهداف استفاده از این داده‌ها است. حجم زیاد این داده‌ها و سردرگمی حاصل از آن مشکلی است که مانع رسیدن به نتایج قابل توجه می‌شود. استفاده از تکنیک‌های داده کاوی در بیماری‌های قلب و عروقی با غلبه بر این مشکل و به دست آوردن روابط مفید بین عوامل خطر زا بسیار مفید بوده است. این بیماری‌ها با توجه به شیوع و سهمی که در مرگ و میر انسان‌ها داشته‌اند از اهمیت بالایی برخوردارند. در ابتدای قرن بیستم ۱۰ درصد کل مرگ و میرها به علت بیماری‌های قلبی عروقی بوده است، که در انتهای همین قرن موارد مرگ و میر ناشی از بیماری‌های قلبی به ۲۵ درصد افزایش یافت و پیش‌بینی می‌شود با توجه به روند

الکتروکاردیوگرام یا نوار قلب: در این روش با استفاده از الکترودهایی که روی قفسه سینه و اندامها قرار می‌گیرد، گزارشی از فعالیت الکتریکی قلب ثبت می‌شود.

اکوکاردیوگرافی: در واقع سونوگرافی مخصوص قلب می‌باشد که با استفاده از امواج مأواه صوت به بررسی ساختمان و کار قلب می‌پردازد. با اکوکاردیوگرافی به سایر حفره‌های قلب، ناهنجاری های مادرزادی قلب، میزان انقباض و بیماری‌های دریچه ای قلب پی می‌بریم. بعد از نوار قلب، اکوکاردیوگرافی معمول‌ترین اقدام تشخیص در بیماری‌های قلب و عروق بوده و کامل کننده اطلاعاتی است که توسط معاینه به دست می‌آید.

تست ورزش: در این روش درحالی که بیمار روی دستگاه دویدن ثابت (تردمیل) حرکت می‌کند از بیمار نوار قلب گرفته می‌شود و در صورت وجود تنگی در عروق قلب، همزمان با افزایش ضربان قلب با ورزش، در نوار قلب بیمار تغییراتی حاصل می‌شود که پزشک را به تشخیص گرفتگی عروق کرونر هدایت می‌کند. در صورت مثبت بودن تست ورزش، بیمار با تشخیص پزشک نیاز به انجام آثیوگرافی عروق کرونر دارد. از فاکتورهای خطر برای بیماری‌های قلبی می‌توان به مواردی نظیر فشار خون بالا، دیابت، کشیدن سیگار، کلسترول بالا، سابقه خانوادگی، چاقی، استفاده از کوکائین و محرک‌های مشابه نام برد.

۳-۲- داده کاوی

داده کاوی استخراج اطلاعات مفهومی، ناشناخته و به صورت بالقوه مفید از پایگاه داده می‌باشد. داده کاوی علم استخراج اطلاعات مفید از پایگاه‌های داده یا مجموعه داده ای می‌باشد [۸]. داده کاوی پل ارتباطی میان علم آمار، علم کامپیوتر، هوش مصنوعی، الگوشناسی، فرآگیری ماشین و بازنمایی بصری داده می‌باشد و فرآیندی پیچیده جهت شناسایی الگوها و مدل‌های صحیح، جدید و به صورت بالقوه مفید، در حجم وسیعی از داده می‌باشد، به طریقی که این الگوها و مدل‌ها برای انسان‌ها قابل درک باشند. داده‌ها اغلب حجمی می‌باشند و به تنها یک قابل استفاده نیستند، بلکه دانش نهفته در داده‌ها قابل استفاده می‌باشد [۹]. بنابراین بهره‌گیری از قدرت فرآیند داده کاوی جهت شناسایی الگوها و مدل‌ها و نیز ارتباط عناصر مختلف در پایگاه داده جهت کشف دانش نهفته در داده‌ها و نهایتاً تبدیل داده به اطلاعات، روز به روز ضروری‌تر می‌شود. تشخیص بیماری‌های مختلف در علم پژوهشی یکی از زمینه‌های پرکاربرد داده کاوی

هستند، نبوده بلکه پیدا کردن عواملی که در بروز این بیماری نقش بیشتری دارند هستیم.

۲- مرور ادبیات

۲-۱- بیماری قلبی و نشانه‌های اولیه آن

اصطلاح بیماری قلبی به تعدادی از بیماری‌هایی که سیستم گردش خون (شامل قلب و رگ‌های خونی) و شیوه‌هایی که در آن خون پمپاژ شده و در سراسر بدن گردش می‌کند و آن را تحت تأثیر قرار می‌دهد اشاره می‌کند و همچنین می‌توان از آن به عنوان حمله قلبی نیز یاد کرد. کاردیومیوپاتی و بیماری‌های قلب و عروق برخی از بیماری‌های قلبی هستند. بیماری قلبی عروقی منجر به بیماری شدید، از کار افتادگی و مرگ می‌شود [۴] و [۵]. پاریک شدن شریان کرونری منجر به کاهش خون و اکسیژن رسانی به قلب و در نتیجه منجر به بیماری کرونری قلب می‌شود. آنفارکتوس میوکارد، عموماً به عنوان حملات قلبی و درد قلبی، و یا درد قفسه سینه شناخته می‌شوند، که توسط بیماری کرونری در بر گرفته می‌شود. انسداد ناگهانی شریان کرونری، عموماً ناشی از لخته شدن خون در یک حمله قلبی است. قفسه سینه وقتی که خون دریافت شده توسط عضلات قلب ناکافی باشد بوجود می‌آید.

فشار خون بالا، بیماری عروق کرونر، بیماری دریچه قلب، سکته مغزی، یا تب روماتیسمی / بیماری روماتیسمی قلب، اشکال مختلف بیماری‌های قلبی و عروقی هستند [۶]. علائم بیماری از مهم‌ترین ابزارهای تشخیص در تعیین نوع بیماری هر فرد می‌باشد و پزشک با تکیه به شرح حال دقیق از روش‌های تشخیصی مکمل استفاده می‌کند.

۲-۲- تشخیص بیماری قلبی

با توجه به نشانه‌های اولیه ای که برای بیماری‌های قلبی وجود دارد پزشکان براساس روش‌های تشخیص که در ادامه معرفی خواهیم نمود بیماری‌های قلبی را تشخیص می‌دهند. از فاکتورهای خطر برای بیماری‌های قلبی می‌توان به مواردی نظیر فشار خون بالا، دیابت، کشیدن سیگار، کلسترول بالا، سابقه خانوادگی، چاقی، استفاده از کوکائین و محرک‌های مشابه نام برد. برخی از روش‌های تشخیص بیماری‌های قلبی را در زیر مشاهده می‌کنید [۷] :

گفته "مايكل ميتزن ماخر" استاد علوم رايانيه اي دانشگاه هاروارد هدف از استفاده از چنین ابزاری بررسی بيشترین ابعاد ممکن در دادهها و بيشترین ارتباطات ممکن در ميان آنها است تا در نهايیت برترین دادهها انتخاب شوند. اين ابزار نوعی ابزار اكتشافي است که می تواند ارتباطات ميان دادهها را کشف کرده و براساس اهمیت رتبه بندی کند. براساس گزارش تی جی دیلی، امكان جستجو را يافتن يك الگو در شرایطی که محقق هنوز نمی دارد دقیقا در جستجوی چه چیزی است، به وجود خواهد آمد و می تواند ایده‌های جدید و ارتباطاتی ایجاد کند که تاکنون کسی به فکر جستجو برای يافتن آنها نیفتاده است. به صورت ویژه این ابزار برای اكتشاف در ميان مجموعه داده‌های مناسب است که می توانند حاوی بيش از يك الگوی مهم باشند. Clementine انواع گوناگونی از روش‌های مدل سازی مانند روش‌های مبتنی بر یادگیری ماشین، هوش مصنوعی و آمار را ارائه می دهد. از طریق این روش‌ها می توان اطلاعات جدیدی را از داده ورودی استخراج نموده و براساس آن مدل‌های پیشگویی مورد نظر را ایجاد کرد. هر روش نقاط قوت خاص خود را دارد و برای حل دسته خاصی از مسائل سودمند است. گره‌های مدل سازی در قالب مازول‌های dataacademy.ir يابه Clustering، Classification و Association تقسیم می شود [۱۳] و [۱۴]. ما در این مقاله بر روی روش کلاسه بندی (Classification) متمرکز شده که اساس کلاسه بندی (Classification) مبتنی بر الگوریتم‌های یادگیری با نظارت، هدف از کاوش داده‌ها است. در الگوریتم‌های یادگیری با نظارت، هدف از کاوش داده‌ها مشخص بوده و می دانیم که به دنبال چه نوع دانشی می گردیم ولی در الگوریتم‌های یادگیری بدون نظارت مانند خوشه بندی (Clustering)، هدف کاملاً تعریف شده نیست.

۳- مجموعه داده ها

مجموعه داده‌های بیماری قلبی از ۴ موقعیت زیر جمع آوری شده اند:

- ✓ کلینیک کلیولند (Cleveland)
- ✓ موسسه مجارستانی قلب و عروق (Hungarian)
- ✓ مرکز پزشکی، لانگ بیچ (long-beach-va)
- ✓ دانشگاه علوم پزشکی زوریخ (Switzerland)

داده‌های تمام این پایگاه داده‌ها فرمت یکسانی دارد. این پایگاه داده‌ها ۷۶ سطر Attribute داشته در حالیکه فقط

محسوب می‌شود. استخراج قواعد طبقه بندی، نوعی داده کاوی است که در آن دانشی به شکل چندین قانون ساده و فهم پذیر از داده کشف شده و در آینده برای تصمیم گیری و پیشگویی به کاربرده می شود [۱۰]. با به کارگیری الگوریتم‌های داده کاوی می توان سیستم‌های هوشمندی ابداع کرد که به شکل خودکار و بدون نیاز به نظارت پزشک قادر به فهم و تفسیر ویژگی‌های پزشکی افراد باشند و یا اطلاعات مفیدی را اکتشاف کنند که متخصصان را در قضایت صحیح یاری رساند [۱۱].

انبار داده به مجموعه ای از داده‌ها گفته می شود که از منابع مختلف اطلاعاتی سازمان جمع آوری، دسته بندی و ذخیره می شود. در واقع یک انبار داده مخزن اصلی کلیه داده های حال و گذشته یک سازمان می باشد که برای همیشه جهت انجام عملیات گزارش گیری و آنالیز در دسترس مدیران می باشد. انبار داده بیماری‌های قلبی شامل غربالگری از اطلاعات مربوط به بیماران قلبی است. در ابتدا، انبار داده‌ها پیش پردازش می شود و در مراحل بعد با استفاده از ابزارهای داده کاوی نظیر کلمانتاین عملیات نهایی بر روی داده‌های انبار داده انجام می شود [۱۲].

۴-۲ Clementine

کلمانتاین یکی از پرکاربردترین ابزارهای داده کاوی است. محقق ایرانی دانشگاه هاروارد به همراه گروهی از محققان در MIT موفق به ابداع این ابزار داده کاوی شده اند که می تواند الگوهای معنا داری را در میان مجموعه ای گسترده از داده‌ها ردیابی کنند. با استفاده از این ابزار می‌توان الگوهای چندگانه پنهان در انواع مجموعه داده‌ها از قبیل داده‌های بهداشتی یا آمار مسابقات بیسبال کشف کرد. ابزار کلمانتاین یک محیط کاری برای انجام فرآیندهای داده کاوی است که به ما کمک می کند مدل‌های پیشگوی را با استفاده از دانش فنی مورد نیاز ساخته و آنها را به منظور تصمیم گیری بهتر در فرآیندهای تجاري به کار بگيريم. این ابزار که براساس مدل استاندارد CRISP-DM طراحی شده است، از تمامی مراحل فرآیند داده کاوی از زمان ورود داده خام تا حاصل شدن نتایج مورد انتظار پشتیبانی می نماید.

محققان این ابزار را بر روی چندین مجموعه بزرگ داده‌ها از جمله، مجموعه داده‌هایی در رابطه با تریلیونها ریزجاذبه‌انی که درون روده زندگی می کنند، آزمایش کردند. این ابزار توانست در حدود ۲۲ میلیون مقایسه را در میان داده‌های دریافتی انجام داده و الگوهای پنهان در آنها را به چند صد الگوی جالب توجه که پیش از این مورد مشاهده قرار نگرفته بودند، محدود سازد. به

۴- الگوریتم های مورد استفاده

۴-۱- ماشین بردار پشتیبان

یکی از روش های یادگیری، یادگیری بانظارت است که از آن برای رده بندی و رگرسیون استفاده می کنند. این روش از جمله روش های نسبتاً جدیدی است که در سال های اخیر کارایی خوبی را نسبت به روش های قدیمی تر برای طبقه بندی از جمله شبکه های عصبی پرسپترون نشان داده است. مبنای کاری دسته بندی کننده ماشین بردار پشتیبان یا SVM دسته بندی خطی داده ها است که در تقسیم خطی داده ها سعی می کنیم خطی را انتخاب کنیم که حاشیه اطمینان بیشتری داشته باشد [۱۵]. حل معادله پیدا کردن خط بهینه برای داده ها به وسیله روش های QP که روش های شناخته شده ای در حل مسائل محدودیت دار هستند صورت می گیرد. قبل از تقسیم خطی برای اینکه ماشین بتواند داده های با پیچیدگی بالا را دسته بندی کند داده ها را به وسیله ای تابع PHI به فضای با ابعاد خیلی بالاتر می برمی.

۴-۲- درخت تصمیم

درخت تصمیم یک ابزار برای پشتیبانی از تصمیم است که از درخت برای مدل سازی استفاده می کند. به طور معمول درخت تصمیم در تحقیق در عملیات استفاده شده و به طور خاص در آنالیز تصمیم و نیز برای تشخیص استراتژی که با بیشترین احتمال به هدف برسد به کاربرده می شود. استفاده دیگر درخت تصمیم، توصیف محاسبات احتمال شرطی است. در میان ابزارهای پشتیبانی تصمیم، درخت تصمیم و دیاگرام تصمیم اداری مزایایی هستند از جمله فهم ساده، کارکردن با داده های بزرگ و پیچیده، استفاده مجدد آسان، در صورتی که درخت تصمیم برای یک مسئله ساخته شد، نمونه های مختلف از آن مسئله را می توان با آن درخت تصمیم محاسبه کرد و همچنین قابلیت ترکیب با روش های دیگر اشاره کرد.

چهار الگوریتم برای ساخت درخت تصمیم موجود هستند. این الگوریتم ها اساساً یک عمل شبیه به هم را پیاده سازی می کنند. آنها تمامی صفات پایگاه داده را بررسی می کنند تا به صفاتی برسند که بهترین رده بندی و پیش بینی را با تقسیم داده به زیر گروه ها انجام می دهد. این فعالیت به صورت بازگشتی تکرار می شود تا باز هم زیر گروه ها به زیر گروه های دیگری شکسته شوند. صفات هدف یا ورودی می توانند از نوع عددی و یا طبقه ای بر حسب الگوریتم مورد استفاده باشند اگر

۱۴ سطر از آنها استفاده می شوند. همه Attribute ها مقدار عددی دارند و مجموعه داده های استفاده شده برای پژوهش ترکیبی از ۴ مجموعه داده بالاست که تعداد آنها ۹۲۰ رکورد می باشد.

مجموعه داده ها شامل Attribute های زیر می باشد :

- Id: شماره شناسایی بیمار

- Age : سن

- Sex : جنسیت (1 = مرد، 0 = زن)

= 0 ، substernal = 1 = محل درد قفسه سینه (1 = otherwise

(1= provoked by exertion; otherwise=0) : Painexer

)1= relieved after rest; otherwise = 0(: Relrest

- Cp : نوع درد قفسه سینه

- Trestbps : فشار خون ایستا

- Chol : کلسترول

- Famhist : سابقه خانوادگی

- Restecg : نتایج کاردیو گرافی

- reading month of exercise: ECG

- Thaldur : مدت زمان تست ورزش

- Thalach : حداکثر ضربان قلب

- Thalrest : ضربان قلب در زمان استراحت

- Num : تشخیص بیماری قلبی (وضعیت آنژیو گرافی)

- Value 0 : احتمال باریک شدن قطر رگ

کمتر از ۵۰٪ است (سالم)

- Value 1 : احتمال باریک شدن قطر رگ

بیشتر از ۵۰٪ است (بیمار)

از ۱۶ صفت ذکر شده صفت ۱ مشخصه شناسایی و صفت ۱۶ به عنوان کلاس است، که در کل از ۱۴ صفت برای پیش بینی استفاده شده است. با توجه به مجموعه داده های بیماری قلب در نمودار ۱ می توان کلاس بیماری قلبی را مشاهده کرد که به دو قسمت سالم (Risky) و بیمار (Healty) تقسیم می شود.

Value /	Proportion	%	Count
healty		44.87	411
risky		55.33	509

نمودار ۱- کلاس بیماری قلب در دو قسمت سالم و بیمار

۴-۳- بیزین

ارائه کننده مدلی پیش بینی کننده در رابطه با احتمال خروجی / نتایج خاصی است. الگوریتم Naive Bayes یا بیزین الگوها و یا ارتباط بین داده ها را با شمارش تعداد وقوع مشاهده آنها اندازه گیری می نماید. این الگوریتم سپس مدلی را ایجاد می کند که نشان دهنده الگوها و ارتباط آنها با یکدیگر است. پس از ایجاد مدل، می توان از آن به عنوان الگوئی پیش بینی کننده استفاده نمود. این الگوریتم کمک خواهد کرد تا بتوان مدل هایی را برای طبقه بندی و پیش بینی اهداف متعددی انجام داد مانند، کدام مشتری ها علاقه مند هستند تا محصول خاصی را خریداری نمایند، کدام مشتری ها قادر هستند بیش از ۳۰۰ دلار خریداری نمایند، شناسائی مشتری هایی که امکان دارد محصولات شرکت را ترک کرده و از محصولات رقبا خریداری نمایند، پیش بینی محصولاتی که میل به شکست داشته و میزان احتمال آن.

الگوریتم Naïve Bayes این پیش بینی ها را توسط تئوری Bayes (که فرض می کند مقادیر داده مستقل از یکدیگر می باشند) انجام می دهد. این الگوریتم ایجاد مدل را در سریع ترین زمانی، ممکن می سازد، که برای مسئله همچون طبقه بندی با کلاس [dataacademy.ir](#) یا بیشتر قابل انجام خواهد بود [۱۷].

۴-۴- لیست تصمیم

با در دسترس داشتن رکوردهای مجموعه یادگیری، این الگوریتم قادر به تولید دو نوع معروف از کلاس بندی های جدول تصمیم است که کار هردو بر اساس یک جدول جستجو می باشد. هر جدول تصمیم دارای ساختاری مشکل از دو قسمت است، الگو؛ شامل لیستی از صفات می باشد و بدنه؛ مجموعه از رکوردهای برچسب دار است. هر رکورد دارای مقادیری برای صفات الگو و یک مقدار برای برچسب می باشد. مجموعه ای رکوردها با مقادیر یکسان برای الگو را، یک سلول می نامیم.

اولین کلاس بند DTMaj نام دارد که در صورت خالی بودن سلولی که با رکورد جدید برابری می کند، کلاس اکثربت داده های یادگیری را بر می گرداند. دومین کلاس بند DTLoc نام دارد که در صورت خالی بودن سلولی که با رکورد جدید برابری می کند، به دنبال نظری بر این صفات کم تر برای آن می گردد. این روش پاسخ را از میان همسایگان محلی بر می گرداند و برای مجموعه داده های یکنواخت کارا می باشد زیرا تغییرات کوچک در صفات، باعث تغییر در برچسب نمی شود. فرآیند

یک بازه مورد استفاده قرار گیرد نتیجه کار یک درخت رگرسیون خواهد بود، اما اگر ورودی ها بصورت رده ای باشند نتیجه کار یک درخت رده بندی خواهد بود [۱۶].

C& R Tree

این درخت یک درخت تصمیم تولید می کند که سعی در پیش بینی و رده بندی مشاهدات آینده دارد. این روش سعی در کم کردن ناخالصی در هر رده دارد. یک گره وقتی کاملاً عاری از ناخالصی است که تمامی عناصر یک زیر گروه آن متعلق به یک رده از فیلد هدف باشند. صفت پیش بینی کننده و فیلد هدف می توانند از دو نوع بازه و رده ای باشند. تمامی تقسیم بندی ها دودویی خواهند بود، به این معنی که فقط دو زیر گروه از هر گره منشعب خواهند شد.

QUEST

این درخت عملکردی همانند درخت C&R دارد، اما سعی دارد زمان لازم برای ساخت درخت تصمیم توسط C&R را کاهش دهد. در ضمن صفت پیش بینی کننده آن مثل C&R می تواند هم طبقه ای و هم بازه باشد، اما فیلد هدف آن باید حتماً رده ای باشد. درخت تولید شده توسط آن یک درخت دودویی است.

CHAID

برخلاف درخت C&R و QUEST، CHAID می تواند درختی تولید کند که در برخی موارد به صورت غیر دودویی عمل کند، یعنی یک گروه آن به سه زیر گروه و یا بیشتر شکسته شود. صفات پیش بینی کننده و هدف می توانند هم از نوع بازه و هم از نوع رده ای باشد.

C5

برای ساخت یک درخت تصمیم و یا مجموعه قوانین استفاده می گردد. این گره بر حسب نیاز می تواند برای برخی گره ها بیش از دو زیر گروه ایجاد کند. صفت هدف آن حتماً باید از نوع رده ای باشد.

در بین ۴ الگوریتم درخت تصمیم ما در ادامه این مقاله از الگوریتم های C&R Tree، C5 و CHAID بهره خواهیم برداشت.

با توجه به مجموعه داده موجود می‌توان از روش‌های مختلفی استفاده کرد ما در این پژوهش برای مدیریت مقادیر از دست رفته از روش استفاده از بیشترین تکرار استفاده کرده ایم بدین صورت که به جای مقادیر از دست رفته اگر داده گسسته باشد از بین مقادیر موجود بیشترین تکرار را به داده از دست رفته و اگر داده پیوسته باشد میانگین مقادیر داده‌های موجود برای آن فیلد را به عنوان مقادیر از دست رفته اختصاص می‌دهیم.

ساخت جدول تصمیم، شامل پیدا کردن یک لیست بهینه از صفات برای الگو می‌باشد. این لیست به گونه‌ای است که کمترین میزان خطای را روی جمعیت رکوردها دارد. برای انتخاب صفات بهینه، از معیار خستگی استفاده می‌شود. می‌توان از معیارهای دیگری نظیر بهره‌ی اطلاعات، نرخ بهره و اندیس جینی نیز استفاده نمود [۱۸].

۴-۵- روش ترکیبی

۴-۵- نتایج آزمایش

در این پژوهش ۷ الگوریتم روش Classification برای پیشگویی صحت(Accuracy) بیماری قلبی مورد استفاده قرار می‌گیرند که عبارتند از SVM، C5، C&R Tree، CHAID، Decision List، Naive Bayes، Tarkib و روش ترکیبی (Tarkib).

بعد از پیش پردازش داده‌ها، داده‌ها را به دو مجموعه داده Train و Test تقسیم می‌کنیم. با استفاده از داده‌های مجموعه Test مدل مورد نظر ما ساخته می‌شود و سپس داده‌های Train خود را بر روی مدل ساخته شده توسط مجموعه داده Train آزمایش می‌کنیم، در این پژوهش با توجه به Attribute dataacademy.ir مجموعه داده را به دو بخش که ۷۰ درصد از داده‌ها برای Training و ۳۰ درصد برای Testing می‌باشد تقسیم می‌شوند.

۷ الگوریتم مطرح شده را بر روی داده‌ها اجرا می‌کنیم، هدف اصلی بدست آوردن Accuracy به وسیله الگوریتم‌های طبقه‌بندی در دو مجموعه داده Train و Test می‌باشد. مجموعه داده‌ها شامل ۹۲۰ نمونه با ۱۴ Attribute مختلف است. این نمونه‌ها در مجموعه داده نشانگر نتیجه‌انواع مختلف تست‌هایی است که برای پیشگویی صحت بیماری قلبی صورت گرفته است می‌باشد. عملکرد کلاسی‌فایرها مورد ارزیابی قرار می‌گیرد و سپس نتایج آن تجزیه و تحلیل می‌شود و عملکرد را می‌توان براساس نرخ خطای مشخص کرد. با مقایسه میان این الگوریتم‌ها به این نتیجه می‌رسیم که با اینکه سعی شده است روش ترکیبی را به صورتی اعمال شود که بالاترین Accuracy را در میان دیگر الگوریتم‌ها داشته باشد ولی بالاترین دقت در بین الگوریتم‌ها برای مجموعه Train مربوط به الگوریتم C5 و برای مجموعه Test مربوط به الگوریتم ترکیبی است. با توجه به نمودار ۲ الگوریتم‌ها C5 و Decision List که ترکیبی از الگوریتم‌های Tarkib، SVM، C5، Decision List و C5، SVM می‌باشد بهتر از الگوریتم‌های دیگر کار می‌کند زیرا نسبت به بقیه صحت بالاتری داشته و الگوریتم Decision List نسبت به

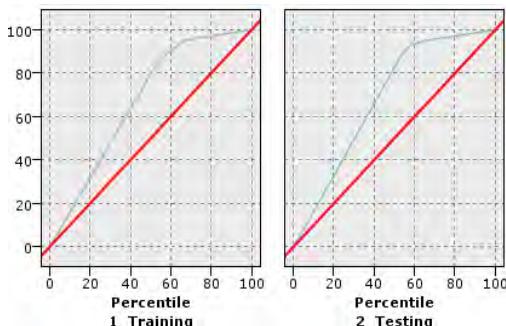
یکی از مشهورترین روش‌های ترکیبی روش رای اکثریت می‌باشد که در بین خروجی‌های کلاسه‌بندها، کلاسی به عنوان نتیجه نهایی انتخاب می‌شود که توسط تعداد کلاسه‌بنده بیشتری به عنوان خروجی پیشنهاد داده باشد. این روش علیرغم سادگی، کارایی تمام کلاسه‌بندها را بدون توجه به ویژگی‌های هر یک، یکسان در نظر می‌گیرد [۱۹]. در این مقاله ما الگوریتم‌های Decision List، SVM، C5 ترکیب کرده ایم تا الگوریتم ترکیبی (Tarkib) بدست آید.

۵- نتایج

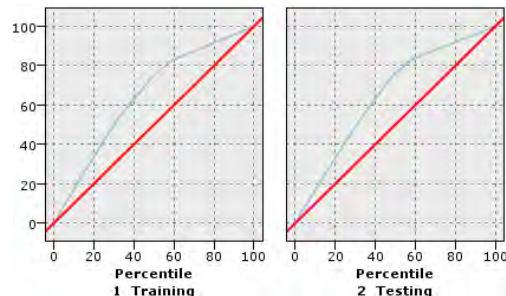
۱-۵- پیش پردازش داده‌ها

پیش پردازش و آماده سازی داده‌ها مهمترین و زمان برترین مرحله در پروژه‌های داده‌کاوی است. تقریباً ۶۰ تا ۹۰ درصد زمان انجام یک پروژه داده‌کاوی صرف این مرحله می‌شود و ۷۵ تا ۹۰ درصد موفقیت پروژه‌ها داده کاوی به آن بستگی دارد [۲۰]. فرآیندهایی که در پیش پردازش انجام می‌شود عبارت است از، تجمعیع، نمونه‌گیری، کاهش بعد، انتخاب زیرمجموعه ویژگی‌ها، ایجاد ویژگی و تبدیلات داده و بر اساس نوع کاربردی که عمل داده کاوی باید بر روی آن انجام شود، تکنیک‌های مختلفی برای هر یک از این اعمال مورد استفاده قرار می‌گیرد [۲۰].

برای پیش پردازش بر روی مجموعه داده بیماری قلبی باید مقادیر از دست رفته یا Missing Value را توسط یکی از روش‌هایی نظیر روش‌های استفاده از بیشترین تکرار، استفاده از میانگین گیری، استفاده از روش نزدیکترین همسایه، حذف کردن کل سطری که داده از دست رفته دارد و یا تکرار فیلدی که داده از دست رفته دارد به ازای مقادیر ممکن، مدیریت کنیم.

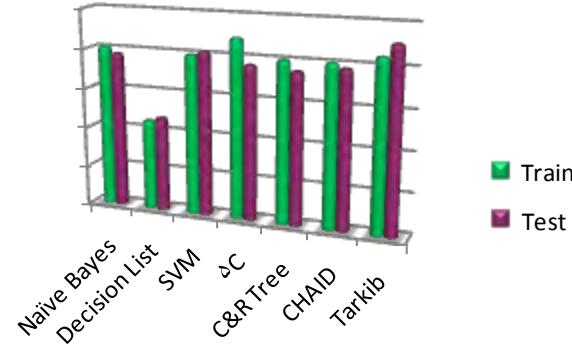


نمودار ۳- منحنی های مربوط به الگوریتم ترکیبی



نمودار ۴- منحنی های مربوط به الگوریتم Decision List

بقیه الگوریتم ها از دقت پایین تری برخوردار می باشد. نمودار ۲ ، مقایسه مقادیر پارامتر میزان درستی مربوط به الگوریتم های اعمال شده روی مجموعه داده بیماری های قلبی را نمایش می دهد. لازم به تأکید است که، از آنجایی که الگوریتم رگرسیون روی این نوع داده ها جواب قابل قبولی نخواهد داد پس این الگوریتم را درنظر نگرفته ایم.



نمودار ۲- مقایسه دقت الگوریتم ها

۴-۵- افزایش دقت

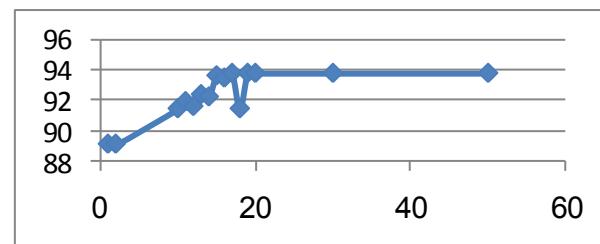
همانطور که در بخش قبل مشاهده شد الگوریتم C5 از دقت بالاتری نسبت به دیگر الگوریتم ها بر روی مجموعه داده Train پیشرفت کرده است. با این بخش به دنبال بالا بردن دقت این برخوردار می باشد و در این بخش به دنبال تکنیک Boosting هستیم. یکی از تکنیک الگوریتم با استفاده از تکنیک Boosting است. یکی از تکنیک هایی که می توان با استفاده از آن میزان دقت را بهبود داد، تکنیک Boosting در درخت تصمیم C5 می باشد، عملکرد این تکنیک بدین صورت است که چندین مدل را در یک دنباله تولید کرده، اولین مدل را با روش معمول C5 می سازد و سپس مدل دوم با تمرکز بر روی رکوردهای غلط رده بندی شده در مدل اول ساخته می شود و مدل سوم با تمرکز بر روی رکوردهای غلط رده بندی شده در مدل دوم ساخته می شود و به همین منوال تا آخر ادامه می یابد.

با توجه به نمودار ۲ الگوریتم C5 دارای دقت ۸۹.۱۳٪ برای داده های Train و دقت ۷۶.۸۱٪ برای داده های Test می باشد که با اعمال تکنیک Boosting با مرتبه تکرار ۱۷ بر روی الگوریتم C5 ، می توان دقت الگوریتم C5 را به ترتیب به ۹۳.۷۹٪ و ۸۲.۶۱٪ افزایش داد و اگر مرتبه تکرار از ۱۷ بالاتر رود از دقت این الگوریتم همانطور که در نمودار ۵ قابل مشاهده است ادامه می یابد.

۳-۵- منحنی یادگیری

منحنی یادگیری یا Learning curve ، نشان دهنده میزان پیشرفت در یادگیری به شکل نمودار از روی امتحان های موفقیت آمیز یادگیرنده است. به عبارت دیگر منحنی یادگیری جنبه ها و قسمت های معینی از پیشرفت در یادگیری است. منحنی یادگیری یکی از راه های ارزیابی الگوریتم های Classification می باشد، هر چه سطح زیر نمودار سنگین تر باشد قدرت یادگیری بیشتر می شود. هر چه شبیه منحنی تندتر باشد و سطح زیر منحنی به یک نزدیکتر شود یادگیری در حد بالایی می باشد و هر چه منحنی به خط قرمز رنگ نزدیکتر شود تصمیم گیری شانسی و بر حسب تصادف خواهد بود. با توجه به نتایج مرحله قبل الگوریتم ترکیبی بالاترین دقت در پیش بینی و الگوریتم Decision List پایین ترین دقت را در بین سایر الگوریتم ها دارا می باشد. از این رو در زیر منحنی های مربوط به هر کدام از الگوریتم ها را برای دو مجموعه داده Train و Test نمایش داده می شوند. در نمودار مربوط به الگوریتم ترکیبی نمودار ۳ مشاهده می شود که هر چه به سمت بالای نمودار می رویم قدرت یادگیری به ۱۰۰ نزدیکتر می شود. نمودار ۴ مربوط به منحنی الگوریتم Decision List می باشد.

- [3] Carlos Ordóñez, "Improving Heart Disease Prediction Using Constrained Association Rules," Seminar Presentation at University of Tokyo, 2004
- [4] Heart Disease from <http://chineseschool.netfirms.com/heart-disease-causes.html>
- [5] "Hospitalization for Heart Attack, Stroke, or Congestive Heart Failure among Persons with Diabetes", Special report: 2001 – 2003, New Mexico.
- [6] Heart disease from http://en.wikipedia.org/wiki/Heart_disease
- [7] Diagnosis Of Heart Disease Using Datamining Algorithm, Asha Rajkumar , Mrs. G.Sophia Reena, Page 38 Vol. 10 Issue 10 Ver. 1.0 September 2010.
- [8] D. Hand, H. Mannila, P. Smyth(2001). Principles of Data Mining. MIT Press, Cambridge.
- [9] Huang T.-M., Kecman V., Kopriva I. (2006), Kernel Based Algorithms for Mining Huge Data Sets, Supervised, Semi-supervised, and Unsupervised Learning, Springer-Verlag, Berlin, Heidelberg, 260 pp. 96 illus., Hardcover, ISBN 3-540-31681-7 .
- [10] Franck Le Duff , Cristian Munteanb, Marc Cuggiaa, Philippe Mabob, "Predicting Survival Causes After Out of Hospital Cardiac Arrest using Data Mining Method", Studies in health technology and informatics, 107(Pt 2):1256-9, 2004.
- [11] Boleslaw Szymanski, Long Han, Mark Embrechts, Alexander Ross, Karsten Sternickel, Lijuan Zhu, "Using Efficient Supanova Kernel For Heart Disease Diagnosis", proc. ANNIE 06, intelligent engineering systems through artificial neural networks, vol. 16, pp.305-310, 2006.
- [12] W. Frawley and G. Piatetsky. Knowledge Discovery I DataBases.ISSN 0738-4602.
- [13] Kiyoung Noh, HeonGyu Lee, Ho-Sun Shon, Bum JuLee, and Keun Ho Ryu, "Associative Classification Approach for Diagnosing Cardiovascular Disease", Springer, Vol:345 , pp721-727, 2006.
- [14] R. Agrawal and R. Srikant, "Fast algorithms for mining association rules in large databases", In Proceedings of the 20th International Conference on Very Large Data Bases, Santiago, Chile, August 29-September 1994.
- [15] Vojislav Kecman: "Learning and Soft Computing Support Vector Machines, Neural Networks, Fuzzy Logic Systems", The MIT Press, Cambridge, MA, 2001.
- [۱۶] مینایی بهروز، نصیری مهدی، حسنی دانیال، شناسا ابراهیم؛ آموزش گام به گام داده کاوی با Clementine؛ انتشارات گروه مهندسی پژوهشی ساحر، انتشارات پارسه، چاپ اول، پاییز ۱۳۹۰.
- [۱۷] جام سحر خشایار، کتاب رهیافتی بر هوشمندی کسب و کار.
- [۱۸] محمد رضا کیوان پور و لیلا خلعتبری، "مقایسه الگوریتم های کلاس بندی در تشخیص دیابت و نارسایی قلبی" مجموعه مقالات سومین کنفرانس داده کاوی ایران، دانشگاه علم و صنعت، آذر ۱۳۸۸.
- [۱۹] حمیدرضا طهماسبی و حسن احمدی، "افزایش دقت کلاس بندی در داده کاوی با استفاده از ترکیب کلاس بندها" مجموعه مقالات سومین کنفرانس داده کاوی ایران، دانشگاه علم و صنعت، آذر ۱۳۸۸.
- [۲۰] علی محمد احمدوند، بهروز مینایی بیدگلی والهام آخرند زاده، "تحلیل رضایت مندی شهروندان با استفاده از تکنیک های داده کاوی" مجموعه مقالات سومین کنفرانس داده کاوی ایران، دانشگاه علم و صنعت، آذر ۱۳۸۸.



نمودار ۵- افزایش دقت الگوریتم C5 با تکنیک Boosting

۶- نتیجه گیری

بررسی ها و مقایسات انجام شده روی مجموعه داده بیماری های قلبی و یافتن بهترین و کارترین و همچنین بدترین الگوریتم ها در تشخیص بیماری ما را به این نتیجه می رسانند که هرگز نمی توان الگوریتمی را به عنوان الگوریتم بهینه معرفی کرد، بر طبق مقاله ای که در سال ۲۰۱۰ در زمینه تشخیص بیماری های قلبی انجام شد [۷] بهترین الگوریتم برای محاسبه بالاترین دقت پیش بینی بیماری های قلبی، الگوریتم Naïve Bayes شناخته شده است در صورتیکه در این پژوهش الگوریتم Naïve Bayes به نسبت دیگر الگوریتم ها دقت پیش بینی پایینی دارد. بنابراین در این بین الگوریتم های ترکیبی جواب قابل قبولتری برای ما ارائه می دهدند.

داده کاوی در مدیریت بهداشت و درمان برخلاف زمینه های دیگر مديون این حقیقت است که داده ها در حال حاضر ناهمگن هستند و اینکه محدودیتهای خاص اخلاقی، حقوقی و اجتماعی درمورد اطلاعات خصوصی پژوهشی وجود دارد. داده های مربوط به مراقبت بهداشت و درمان به طور طبیعی حجمی هستند و آنها از منابع گوناگون بدست می آیند و همه آنها به طور کامل در ساختار یا کیفیت مناسب نیستند. امروزه، بهره برداری از داده های غربالگری بالینی بیماران تجربه متخصصان متعدد و داده های غربالگری بالینی بیماران جمع آوری شده در یک پایگاه داده در طی روند تشخیص، به طور گستردگی استفاده می شود.

مراجع

- [۱] ماهنامه تخصصی پژوهشی شماره ۱۳۹۱، اسفند ماه ۱۳۹۰
- [۲] M. Chen, J. Han, and P. S. Yu. Data Mining: An Overview from Database Perspective. IEEE Trans. Knowl. Dat. Eng., vol: 8, no:6, pp: 866-883, 1996.