

## طراحی سیستم پیش‌بینی بیماری قلبی-عروقی با استفاده از الگوریتم‌های یادگیری ماشین

محمد کاظمی<sup>۱</sup>، امیرعباس کاظمی<sup>۲</sup>، امیرحسین باب الحکمی<sup>۳</sup>، نگار عباسی<sup>۴</sup>

<sup>۱</sup>دانشجوی کارشناسی ارشد، گروه کامپیوتر، دانشگاه آزاد اسلامی مشهد

<sup>۲</sup>دانشجوی کارشناسی ارشد، گروه کامپیوتر، دانشگاه فردوسی مشهد

<sup>۳</sup>دانشجوی دکتری، گروه کامپیوتر دانشگاه فردوسی مشهد

<sup>۴</sup>کارشناسی ارشد، گروه آمار دانشگاه پیام نور مشهد

### چکیده

بیماری‌های قلبی یکی از شایع‌ترین بیماری‌های است که در حال حاضر تعداد افراد مبتلا به این نوع بیماری‌ها در حال افزایش می‌باشد. این در حالی است در صورتی که مراقبت‌های لازم برای بیمار در زمان مناسب صورت نگیرد، می‌تواند باعث مرگ بیمار شود. از این‌رو تشخیص دقیق در مرحله معاینه اولیه به همراه درمان مناسب می‌تواند منجر به اجتناب از افزایش میزان مرگ‌ومیر ناشی از بیماری قلبی گردد. برای رسیدن به این مهم می‌توان از تکنیک‌های موجود در زمینه داده‌کاوی بهره گرفت. داده‌کاوی داده‌های مفیدی را از مجموعه داده‌های موجود استخراج می‌کند که منجر به پیش‌بینی یا دسته‌بندی اطلاعات از طریق خوشبندی، کلاس‌بندی و یا کشف الگوهای پنهان می‌شود. تاکنون تحقیقات زیادی با استفاده از مدل‌های داده‌کاوی در تشخیص بیماری‌های مختلف مانند بیماری‌های قلبی و عروقی انجام شده است..

در این مقاله قصد داریم با استفاده از رویکردی مبتنی بر انتخاب ویژگی به عنوان یک گام پیش پردازش، مدلی باهدف تشخیص بیماری قلبی ارائه گردد. راهکار پیشنهادی دارای ۳ گام اصلی می‌باشد که گام ۱) پیش پردازش داده‌ها با هدف رفع مقدادر Null و پرت در مجموعه داده‌ها، گام ۲) انتخاب ویژگی‌های موثر با بهره وری از ۲ روش ضریب همبستگی پیرسون و تجزیه و تحلیل مولفه‌های اصلی که سعی در حذف ویژگی‌هایی که با صفت هدف رابطه خاصی ندارند و رفتار این ویژگی مستقل از صفت هدف می‌باشند، است. و در گام ۳) با استفاده از ۳ الگوریتم J48، SVM شبکه بیزین و ۴۸ J با دقت ۰،۸۹٪ دارای بالاترین دقت است.

**واژه‌های کلیدی:** تشخیص بیماری قلبی، انتخاب ویژگی، ضریب همبستگی پیرسون، تحلیل مولفه‌های اساسی، طبقه‌بندی

## ۱. مقدمه

بیماری‌های قلبی یکی از شایع‌ترین بیماری‌ها می‌باشد که در حال حاضر تعداد افراد مبتلا به بیماران قلبی در حال افزایش است. به طور کلی این بیماری دارای علائم مختلفی می‌باشد اما از مهمترین آن می‌توان به درد در قفسه سینه، ناحیه فک، گردن، نفس تنگی و ... اشاره کرد. همچنین مشکلات قلبی می‌تواند باعث بروز بیماری عروق کرونر قلب، نارسایی قلبی و سکته مغزی شود. با این وجود، اگرچه عامل اصلی مرگ و میر در دهه اخیر در دنیا شناخته شده است، به عنوان کنترل پذیرترین و قابل پیشگیری ترین بیماری نیز شناخته می‌شود. نتایج بدست آمده نشان می‌دهد، شناسایی و تشخیص به موقع بیماری قلبی علاوه بر امکان جلوگیری از پیشرفت و بروز عوارض ناشی از آن، می‌تواند در درمان صحیح و کامل بیماری موثر باشد. برای تشخیص به موقع بیماری قلبی تست‌های مختلفی مانند نوار قلب، تست استرس و آنژیو قلب وجود دارد ولی این نوع آزمایش‌ها هزینه‌ی زیادی در بر دارند و نمی‌توان آنها را به طور گستردۀ استفاده کرد. تشخیص در مراحل ابتدایی کمک شایانی به روند درمان افراد میکند و نسبت به روش‌های متدالو پزشکی مفروض به صرفه‌تر است. امروزه داده کاوی و کشف دانش به کمک ما آمده است تا با جمع آوری اطلاعاتی مرتبط، به پیش‌بینی بیماری قلبی بپردازیم. یکی از عملکردهای پیش‌گویانه در داده کاوی، دسته بندی است. دسته بندی، فرآیند یافتن مدلی است که با تشخیص دسته‌ها و یا مفاهیم داده می‌تواند دسته ناشناخته اشیا دیگر را پیش‌گویی کند<sup>[۱]</sup>. هدف از پیش‌بینی، کمک به روند تشخیص و بهبود سلامت بیماران است. تاکنون مطالعات زیادی پیرامون تکنیک‌های داده کاوی بر روی انباره بیماران قلبی و عروقی انجام شده است، که تکنیک‌های مختلف داده کاوی مانند درخت تصمیم و شبکه‌های عصبی را به کار برده اند که (a) پیش‌بینی حمله قلبی و سندروم کرونری حاد با استفاده از شبکه عصبی داده کاوی<sup>[۴]-[۲]</sup> (b) ارزیابی معیارهای مهم برای پیش‌بینی بقا بیمار با استفاده از شبکه‌های بیزین<sup>[۵]</sup> (c) تشخیص ایسکمی قلبی در رکوردهای طولانی مدت الکتروکاردیوگرام<sup>[۶]</sup> (d) ارزیابی فاکتورهای ریسکی مرتبط با بیماری قلبی با استفاده از درخت تصمیم و قوانین تلازمی<sup>[۷]-[۸]</sup> و (e) کشف روابط بین عوامل خطرزای قلبی با استفاده از درخت تصمیم<sup>[۹]</sup> از جمله پژوهش‌های انجام شده در این زمینه است.

در این مقاله نیز قصد داریم با بهره وری از تکنیک‌های داده کاوی و تحلیل ویژگی‌های مختلف مدلی باهدف پیش‌بینی بیماری قلبی ارائه شود. گام‌های ارائه شده شامل ۲ فاز اصلی می‌باشند. در ابتدا (۱) پیش‌پردازش‌های لازم باهدف رفع مقادیر Null انجام گرفته است. برای رفع مقادیر Null از تکنیک پیش‌بینی و طبقه‌بندی نمونه‌ها استفاده شده است. به این صورت که نمونه‌هایی که مقادیر Null دارند را به عنوان داده‌های تست و سایر نمونه‌ها را به عنوان داده‌های آموزش در نظر گرفته‌ایم. (۲) در گام دوم به بررسی ویژگی‌های موثر در پیش‌بینی بیماری قلبی پرداخته خواهد شد و با بهره وری از ۲ روش ضربی همبستگی پرسون و PCA صفاتی که باهدف مسئله ارتباطی ندارند، شناسایی و حذف خواهند شد. (۳) و در گام نهایی از ۳ الگوریتم درخت تصمیم، شبکه بیزین و SVM باهدف ساخت مدل طبقه‌بندی، پیش‌بینی بیماری قلبی استفاده خواهد شد. در ادامه این مقاله، در بخش ۲ مقالات مرتبط به هدف مسئله مورد بررسی قرار خواهند گرفت. در بخش ۳ راهکار پیشنهادی و گام‌های ارائه شده به صورت دقیق‌تر مورد بررسی قرار خواهند شد. در بخش ۴ نتایج آزمایشات و ارزیابی‌های انجام شده ارائه و در انتها نیز نتایج بدست آمده تحلیل و کارهای آینده معرفی می‌شوند.

## ۲. پیشینه

در این بخش از مقاله، برخی از پژوهش‌های مرتبط در پیش‌بینی بیماری قلبی مورد بررسی قرار خواهد گرفت. قلب یکی از مهمترین اجزاء بدن می‌باشد که وظیفه انتشار خوب به سایر بخش‌های بدن را بعده دارد. بافت قلب از ۳ لایه داخلی (آندوکارد)، لایه عضلانی میانی (میوکارد) و لایه خارجی (پریکارد) تشکیل شده است و هر یک از این لایه وظایف مخصوص به خود را دارد. لایه اندکارد پوشش حفره‌های دهلیز و بطن، میوکارد ضخیم‌ترین لایه قلب است که ادامه حیات انسان به آن وابسته می‌باشد و پریکارد پوشش پیوندی یا آبسامه قلب را می‌سازد<sup>[۱]</sup>. در یک تقسیم بندی کلی، بیماری قلبی و عروقی

شامل بیماری عروق کرونری، بیماری دریچه قلب، بیماری های مادرزادی قلب، بیماری عضله قلب و ... می باشد که از این بین بیماری شرایین کرونری از شایعترین بیماری قلبی است. این بیماری عمدتاً به صورت آثین صدری و انفارکتوس قلبی است که از مهمترین دلایل بروز آن می توان به آتروسکلروزیس، تروموبوزیس، اسپاسم و انوریسم کرونری اشاره کرد. آتروسکلروزیس که به معنی اختلال عمومی در شریان هاست، از شایعترین علل رخداد این بیماری می باشد. این اختلال به وسیله پلاک های زردرنگ از کلسترول و چربی ها مشخص می شود و مانع از رساندن اکسیژن به باقت می شود. به عبارتی عروق کرونری تغذیه کننده بافت های قلبی می باشند. وسعت این بیماری به آناتومی توزیع رگ های مسدود بستگی دارد که به عبارتی عروق کرونری تغذیه کننده بافت های قلبی می باشند. با توجه به پیچیدگی های بسیاری که در تشخیص وجود بیماری قلبی یا نوع آن دیده می شود، وجود یک ابزار کمکی که بتواند دانش ارزشمند را در اختیار پزشکان قرار بدهد بسیار احساس می شود تا با بهره وری از این ابزار، بتوان تاثیر داروهای ناسازگار را کاهش و در پیش بینی رفتار آینده بیماران براساس سابقه ثبت شده و تشخیص بیماری قبلی را استفاده کرد.

داده کاوی یک ابزار کاربردی در شناسایی اطلاعات مفید و ناشناخته از حجم عظیم داده هایت. از این اطلاعات می توان برای پیش بینی آینده موقعیت‌ها و به عنوان یک کمک برای فرآیند تصمیم گیری استفاده کرد. دانش مفید را می توان با استفاده از داده کاوی در برنامه کاربردی مراقبت پزشکی مثل سیستم پشتیبانی تصمیم، پیدا نمود. داده بزرگ تحويل داده شده توسط سازمان‌های مراقبت پزشکی بسیار پیچیده و بسیار بزرگ هستند و نمی‌توان آنها را با تکنیک‌های معمول مدیریت و تحلیل نمود. داده کاوی روشی برای تغییر این مقدار داده به اطلاعات مفید برای پشتیبانی تصمیم ارائه می‌کند. داده بزرگ کاوی در مراقبت پزشکی درباره یادگیری مدل‌هایی برای پیش بینی بیماری بیمار است. به عنوان نمونه، داده کاوی می‌تواند به سازمان‌های بیمه سلامت در تشخیص سوء استفاده‌ها یا تظاهر، به مؤسسات پزشکی در تصمیم گیری در مورد مدیریت روابط مشتریان، به پزشکان در دریافت خدمات پزشکی بهتر و اقتصادی‌تر، کمک نماید. این تحلیل پیشگویانه در خدمات پزشکی بسیار استفاده می‌شود. در این مقاله نیز قصد داریم مدلی باهدف پیش بینی بیماری قلبی ارائه شود. اما در این فرآیند از چند تکنیک فیچرسلکشن باهدف آنالیز پارامترهای مختلف استفاده شده است. که به وسیله آن بتوان ویژگی های نامرتب را قبل از آنکه تاثیر منفی بر روی مدل پیش بینی بگذارد شناسایی و حذف کنیم. همچنین براساس خروجی این روش ها از میزان تاثیر هریک از پارامترهای مختلف آگاهی پیدا کرد. با این وجود با توجه به اهمیت این بیماری، مطالعات مختلفی در این زمینه انجام شده است و در ادامه به مرور برخی از این مقالات پرداخته خواهد شد.

در مطالعه ای برای پیش بینی بیماری قلبی از آزمون اکوکاردیوگرافی قفسه سینه استفاده می‌شود که این داده ها در بازه زمانی ۲ سال تهیه شده است. داده ها شامل ۱۵ متفاوت است. چهار آزمایش برای الگوریتم های طبقه بندی طراحی شد و آزمایشها بر روی مجموعه داده های آموزشی کامل شامل ۷۳۳۹ نمونه انجام شد. برای کلیه آزمایش ها دو موقعیت در نظر گرفته شد که حاوی ۱۵ ویژگی و دیگری شامل ۸ ویژگی انتخاب شده است. درمجموع هشت مدل توسعه داده شد. به عنوان نتیجه طبقه بندی J<sub>۴۸</sub> را بر روی مجموعه داده با ۸ ویژگی انتخاب شده با دقت ۹۵,۵۶٪<sup>۳</sup> بعنوان بهترین الگوریتم دسته بندی معرفی نمودند[۲]. در پژوهشی دیگر از الگوریتم های J<sub>۴۸</sub>، شبکه بیزین و بیزین ساده و الگوریتم Simple CART برای طبقه بندی و توسعه یک مدل برای تشخیص حملات قلبی در مجموعه داده بیماران استفاده شد. هدف این پژوهش، پیش‌بینی احتمال حمله قلبی در مجموعه داده بیمار با استفاده از تکنیک های داده کاوی و تعیین مدلی که بالاترین درصد صحت پیش‌بینی را به همراه دارد. مجموعه داده بیماران از داده های پزشکان در آفریقای جنوبی جمع آوری و از ۱۱ ویژگی آن در آزمایشات استفاده گردید[۴]. مقاله ای در سال ۲۰۱۴ با عنوان پیش بینی بستری شدن بدلیل بیماری قلبی با استفاده از روش یادگیری نظارت شده ارائه گردید. داده های مورد استفاده در آن پژوهش از مرکز پزشکی بوستون در بازه زمانی ۵ سال انتخاب شدند که به عنوان نتیجه اعلام گردید که ۸۲٪ بیماران مورد مطالعه در سال بعدی در بیمارستان بستری می شوند[۳] در حوزه های دیگر و با استفاده از الگوریتم ژنتیک و شبکه RBF در سال ۲۰۱۵ پژوهشی صورت پذیرفت که با استفاده از پایگاه داده موجود در سایت معتبر UCI و با ۱۴ ویژگی و ۳۰۳ نمونه است. الگوریتم های طبقه بندی روی داده ها اعمال و به عنوان نتیجه ارائه شد که نتایج

خوبی از دو الگوریتم Naïve Bayes با ۸۳٪ و J48 با ۷۷٪ دقت دریافت شده است [۵]. در مقاله [۶]، مدلی برای توسعه یک شبکه عصبی مصنوعی برای مدل های تشخیص بیماری های عروق کرونری قلب، با استفاده از مجموعه ای از عوامل سنتی و ژنتیکی این بیماری است. پایگاه داده اصلی برای شبکه های عصبی مصنوعی شامل بالینی، آزمایشگاهی، کاربردی، آژیوگرافی عروق کرونر و ژنتیکی است. بهترین دقت با توبولوژی شبکه های عصبی پرسپترون چند لایه از دو لایه پنهان برای مدل شامل شده توسط هر دو عوامل خطر CHD ۱۲ ژنتیکی و غیر ژنتیکی به دست آمد. مجموعه داده استفاده شده در این مقاله، شامل ۹۸۷ بیمار است و تحت آژیوگرافی برای تشخیص عروق کرونر قرار گرفته اند. تشخیص بیماری عروق کرونر قلب نیز از هر دو نتایج بالینی و کرونر گرافی کسب گردید و این نتایج با استفاده از آزمایش و نوع ژنتیک اجازه ایجاد یک پایگاه داده از بیماران مبتلا را فراهم می کند که پس از آن برای تشخیص CHD با استفاده از شبکه های عصبی مورد استفاده قرار بگیرد. مدل شبکه عصبی با استفاده از پرسپترون چند لایه ایجاد و دقت مدل های الگوریتم ژنتیک بهبود یافته شده با استفاده از پارامترهای مختلف بهینه سازی از جمله تعداد نرون در لایه پنهان، تعداد ورودی به شبکه های عصبی و ضریب شیب فعالسازی مورد استفاده قرار گرفته که امکان بهینه سازی را ایجاد می کند. در مقاله [۸] از یک الگوریتم یادگیری ترکیبی برای شبکه عصبی پرسپترون چند لایه فازی از ترکیب دو الگوریتم فرآیندکاری جستجوی گرانشی و بهینه سازی از دحام ذرات استفاده می کنند. ۵ مجموعه داده مربوط به بیماری قلبی، سرطان سینه، هپاتیت، اختلالات کبدی و سرطان ریه در این مطالعه بررسی شده اند. که دقت بدست آمده از اجرای راهکار پیشنهادی بر روی مجموعه داده بیماری قلبی برابر ۷۹,۹۶٪ است.

### ۳. راهکار پیشنهادی

راهکار پیشنهادی ارائه شده برای این پژوهش در ۳ گام تعریف شده است: (۱) پیش پردازش داده‌ها (۲) تعیین مجموعه ویژگی‌های مؤثر و شناسایی شده با ۲ روش پیرسون و PCA و (۳) ارائه مدل پیش بینی که از ۳ الگوریتم درخت تصمیم، شبکه بیزین و SVM استفاده شده است. در ادامه گام ۲ که باهدف انتخاب ویژگی‌های مؤثر و کاهش ابعاد دیتا است و گام ۳ که شامل طبقه‌بندی و آموزش داده‌ها است به صورت دقیق تر مورد بررسی قرار خواهد گرفت.

#### ۳,۱- ضریب همبستگی پیرسون

هدف از استفاده از این روش، شناسایی صفات کم اهمیت برای مسئله هدف و شناسایی مواردی که با یکدیگر همپوشانی دارند است [۱۰]. برای کسب این اطلاعات لازم است ۲ گام انجام شود که شامل:

- ضریب همبستگی صفات با صفت هدف: اگر ضریب همبستگی بین ۰,۰ تا ۰,۲ باشد، یعنی ویژگی مستقل است و تاثیری در مسئله هدف ندارد و حذف می‌شود.

- ضریب همبستگی دو به دوی صفات (جز صفت هدف): اگر ضریب همبستگی بین ۰,۸ و ۰,۲ ویژگی بیش از ۰,۸ باشد یعنی با یکدیگر همپوشانی دارند و می‌توان یکی از صفات را حذف کرد (صفت با ضریب همبستگی کمتر با صفت هدف حذف خواهد شد).

براساس فازهایی که تعیین شد، مجموعه ویژگی‌های منتخب توسط این روش باقیمانده با صفت هدف ارتباط مستقیم (ضریب همبستگی مثبت) و یا ارتباط معکوس (ضریب همبستگی منفی) دارند. همچنین این ویژگی‌ها نسبت به یکدیگر مستقل می‌باشند.

#### ۳,۲- تحلیل مولفه‌های اصلی

دومین روش Feature Selection در گام انتخاب ویژگی‌ها، الگوریتم PCA که یکی از روش‌های پرکاربرد در کاهش ابعاد داده ورودی است، استفاده خواهد شد. تحلیل مولفه‌های اصلی در تعریف ریاضی یک تبدیل خطی متعامد است که داده را به دستگاه مختصات جدید می‌برد به طوری که بزرگترین واریانس داده بر روی اولین محور مختصات، دومین بزرگترین واریانس

بر روی دومین محور مختصات قرار می‌گیرد و همین طور برای بقیه تحلیل مولفه‌های اصلی می‌تواند برای کاهش ابعاد داده مورد استفاده قرار بگیرد، به این ترتیب، مولفه‌هایی از مجموعه داده را که بیشترین تأثیر در واریانس را دارند حفظ می‌کند.

### ۳.۳- الگوریتم‌های طبقه‌بندی

بعد از شناسایی ویژگی‌های مؤثر و حذف ویژگی‌های کم اهمیت، در انتهای داده‌ها را برای طبقه‌بندی و ساخت مدل پیش‌بینی با استفاده از روش طبقه‌بندی درخت تصمیم، شبکه بیزین و SVM آموزش خواهیم داد و براساس نتایج بدست آمده، بهترین مدل از نظر دقت پیش‌بینی به عنوان مدل نهایی ارائه و مشخص خواهد گردید. در ادامه این چند الگوریتم به صورت خلاصه معرفی خواهند شد.

#### ۳.۳.۱- درخت تصمیم

درخت تصمیم یکی از روش‌های قوی و متداول برای دسته بندی و پیش‌بینی است. در واقع درخت‌های تصمیم بالا به پایین رایج‌ترین تکنیک دسته بندی هستند و از مهم‌ترین دلایل رایج بودنشان می‌توان به شفاف بودن، قابل فهم، انعطاف پذیری و پردازش نسبتاً سریع ساختار آنها نام برد. پیش‌بینی به دست آمده از درخت در قالب یک سری قواعد توضیح داده می‌شود. در این درخت هر گره داخلی شامل سوالی بر مبنای یک صفت مشخص و یک فرزند برای هر پاسخ ممکن بوده و هر برگ با یکی از کلاس‌های ممکن برچسب گذاری می‌شود [۱۱]

#### ۳.۳.۲- شبکه بیزین

شبکه بیزین یک مدل گرافیکی برای روابط بین مجموعه‌ای از انواع پارامترهای متغیر است. این مدل گرافیکی ساختار S یک گراف بدون دور جهت دار است و همه گره‌ها در S تناظر یک به یک با ویژگی‌های مجموعه داده دارند. کمان‌ها نشان دهنده تاثیرات بین ویژگی‌ها هستند. زمانی که دسته بندی بیزین برای مجموعه داده‌های بزرگ به کار برده شود دارای سرعت و دقت بالایی است. طبقه‌بندی بیزین ساده نیز بر اساس نظریه بیز است. این روش برای مقادیر ویژگی‌های پیوسته ورودی و خروجی مناسب است و در دنیای واقعی برای تشخیص دست خط کاربرد دارد. این روش مناسب‌ترین روش برای انتخاب یک مدل در مقایسه با داده‌های موجود بر اساس شرایط احتمالی است [۱۲]

#### ۳.۳.۳- ماشین بردار پشتیبان

به مجموعه‌ای از نقاط در فضای  $n$  بعدی داده‌ها، بردار پشتیبان گفته می‌شود که مرز بندی دسته‌ها را نشان داده و دسته بندی و مرزبندی آنها را انجام می‌دهد و با جابجایی یکی از این دو مورد ممکن است تغییر کند. SVM یا ماشین بردار پشتیبان، با معیار قرار دادن بردار‌های پشتیبان بهترین دسته بندی و تفکیک بین داده‌ها را انجام می‌دهد همچنین در SVM مبنای یادگیری ماشین و ساخت مدل، داده‌های قرار گرفته شده در بردارهای پشتیبان می‌باشد. هدف الگوریتم SVM یافتن بهترین مرز در بین داده‌ها بوده و بیشترین فاصله ممکن از تمام دسته‌ها را در نظر می‌گیرد و به سایر نقاط داده‌ها حساس نمی‌باشد.

### ۴- نتایج آزمایشات

در این بخش از مقاله نتایج بدست آمده از اجرای راهکار پیشنهادی بر روی دیتاست ارائه و مورد بررسی قرار خواهد گرفت. برای پیاده‌سازی عملیات پیش‌پردازش، Feature Selection، طبقه‌بندی و ساخت مدل پیش‌بینی از زبان برنامه نویسی R که یکی از زبان‌های پرکاربرد در زمینه داده‌کاوی می‌باشد، استفاده گردید.

#### ۴.۱- مجموعه داده‌ها

داده‌های مورد بررسی در این پژوهش از مجموعه پایگاه داده در این مقاله از مجموعه داده تشخیص بیماری قلبی Cleveland که در سایت مرجع UCI قرار گرفته، استفاده شده است. علائم زیادی از بیماری قلبی وجود دارد و یافتن الگوهایی از داده بیماری قلبی در تشخیص دلایل آتی این بیماری کمک می‌کند [۱۳]. پایگاه داده شامل ۳۰۳ نمونه که ۶ مورد آن دارای مقادیر Null است. این مجموعه داده‌های اولیه ۷۶ صفت خام دارد در حالی که همه آزمایش‌ها فقط بر روی ۱۴ صفت از آن‌ها انجام شده

است. بنابراین، این پایگاه داده شامل ۱۳ علامت بیماری و یک صفت تشخیص است که صفت هدف به وجود بیماری قلبی بر اساس علائم موجود در بیمار اشاره دارد که یک مقدار عددی (کمتر از ۵۰٪ تنگی عروق) یا ۱ (بیشتر از ۵۰٪ تنگی عروق) است.

#### ۴.۲- معیارهای ارزیابی

پس از اجرای مدل سازی باید به ارزیابی نتایج حاصل پرداخت. نتایج ارزیابی باعث بهبود مدل می‌شود و مدل را قابل استفاده می‌نماید. شاخص‌های مختلفی مانند شفافیت (Specificity)، حساسیت (Sensitivity)، دقت (Precision) و صحت (Accuracy) برای ارزیابی روش‌های دسته بندی وجود دارند که طبق روابط ۱ تا ۴ محاسبه می‌گردند. برای محاسبه میزان شاخص‌ها می‌توان از ماتریس اغتشاش (Confusion Matrix) استفاده کرد. این ماتریس ابزار مفیدی برای تحلیل چگونگی عملکرد روش دسته بندی در تشخیص داده‌ها با مشاهدات دسته‌های مختلف است. حالت ایده آل این است که بیشتر داده‌های مرتبط با مشاهدات روی قطر اصلی ماتریس قرار گرفته باشند و مابقی مقادیر ماتریس صفر یا نزدیک صفر باشند. همچنین برای ارزیابی دقیق کارایی از تکنیک K-fold cross Validation استفاده شده است.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (1)$$

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}} \quad (2)$$

$$\text{Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (3)$$

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}} \quad (4)$$

#### ۴.۳- نتایج ضریب همبستگی پیرسون

یکی از مهمترین گام‌های ارائه شده در پیش‌پردازش داده‌ها، آنالیز و تعیین ویژگی‌های موثر می‌باشد. از مهمترین اهداف مدنظر در این گام، حذف ویژگی‌های غیرمرتبط با هدف مسئله که تشخیص بیماری قلبی است، می‌باشد. وجود ویژگی‌های نامناسب و مخرب می‌تواند در ساخت مدل پیش‌بینی تاثیر منفی داشته باشد و دقت مدل را کاهش دهد. همچنین بررسی این پارامترها می‌تواند اطلاعات مفیدی از نحوه ارتباط این پارامترها با یکدیگر و میزان اثرگذاری بر ویژگی هدف را در اختیار پزشکان قرار بدهد. همانطور که در راهکار پیشنهادی ارائه گردید برای این بخش از ۲ الگوریتم پیرسون و PCA استفاده شده است. در جدول ۱ نتایج بدست آمده از اجرای مجموعه داده کلیولند با تکنیک‌های معروفی شده می‌باشد که با توجه به مقدار ضریبی دریافتی رتبه بندی شده است. با توجه به رتبه‌های تشخیص داده شده توسط این ۲ الگوریتم، ۵ ویژگی Trestbps، Restecg، Chol، Fbs و Gender بصورت مشترک به عنوان ویژگی‌های نامطلوب انتخاب و از لیست پارامترهای موجود حذف گردید.

جدول ۱: نتایج حاصل از آنالیز ویژگی‌های مختلف در تشخیص بیماری قلبی

Name	Pearson (Rank)	PCA (Rank)	میزان تاثیر
<b>Age</b>	۰,۲۸ (۸)	۰,۲۸ (۸)	Weak
<b>Gender</b>	۰,۱۹ (۹)	۰,۱۱ (۱۱)	Weak
<b>CP</b>	۰,۴ (۶)	۰,۲۹ (۷)	Moderate
<b>Trestbps</b>	۰,۱۵ (۱۱)	۰,۱۶ (۹)	Independence
<b>Chol</b>	۰,۰۸ (۱۲)	۰,۰۹ (۱۲)	Independence
<b>Fbs</b>	۰,۰۲ (۱۳)	۰,۰۸ (۱۳)	Independence
<b>Restecg</b>	۰,۱۶ (۱۰)	۰,۱۵ (۱۰)	Independence
<b>Thalach</b>	۰,۴۶ (۵)	۰,۳۹ (۲)	Strong
<b>Exang</b>	۰,۴۳ (۳)	۰,۳۳ (۵)	Strong
<b>Oldpeak</b>	۰,۴۲ (۴)	۰,۴ (۱)	Strong
<b>Slope</b>	۰,۳۳۹۲ (۷)	۰,۳۵ (۳)	Moderate
<b>Ca</b>	۰,۴۶ (۲)	۰,۳۱ (۶)	Strong
<b>Thal</b>	۰,۵۳ (۱)	۰,۳۴ (۴)	Strong

#### ۴.۴- نتایج مدل طبقه‌بندی

بعد از شناسایی ویژگی‌های مؤثر، حذف ویژگی‌های کم اهمیت و خوش‌بندی داده‌ها در گام نهایی با استفاده از الگوریتم‌های طبقه‌بندی J48 شبکه بیزین، SVM مدلی باهدف تشخیص وجود بیماری قلبی ساخته خواهد شد و براساس نتایج بدست آمده، بهترین مدل از نظر دقیقت پیش‌بینی به عنوان مدل نهایی ارائه و مشخص خواهد گردید. مجموعه داده ورودی برای این گام شامل ۳۰۳ نمونه و ۹ ویژگی (۸ ویژگی موثر و صفت هدف) است.

باتوجه به دقیقت آمده از آموزش داده‌های بیماری قلبی-عروقی نشان متوسط مدل‌های پیش‌بینی، الگوریتم J48 که یکی از روش‌های درخت تصمیم است، دارای بالاترین دقیقت که برابر ۰,۸۹ نمونه‌های دیتاست را به درستی پیش‌بینی کرده است.

جدول ۲: نتایج تکمیلی الگوریتم‌های طبقه‌بندی

	صفت	دقیقت	حساسیت	شفافیت
<b>J48</b>	۰,۸۹	۰,۹۴	۰,۹۲	۰,۹۲
<b>Bayesian Network</b>	۰,۷۹	۰,۸۸	۰,۸۵	۰,۸۶
<b>SVM</b>	۰,۸۳	۰,۸۸	۰,۸۶	۰,۸۵

## ۵. نتیجه‌گیری

محیط مراکز بهداشتی و درمانی به تعداد زیاد داده در رابطه با بیماران قلبی و عروقی دارند اما یک فقدان ابزار تحلیل موثر برای کشف رابطه‌های مخفی بین بیماران قلبی و عروقی وجود دارد. پس داده کاوی می‌تواند به عنوان ابزار موثری در پیدا کردن اطلاعات پنهان بین بیماران باشد. تکنیک‌های داده کاوی برای پزشکان این کار را ممکن می‌سازد که اطلاعات تشخیص انواع بیماران با شرایط یکسان را پیش بینی نماید. از این رو در این مقاله به ارائه روشی به منظور پیش بینی بیماری قلبی و عروقی پرداخته شده است. گام‌های ارائه شده شامل ۳ فاز اصلی می‌باشند. در ابتدا پیش پردازش‌های لازم باهدف رفع مقدایر Null انجام گرفته است. برای رفع مقدایر Null از تکنیک پیش بینی و طبقه‌بندی نمونه‌ها استفاده شده است. به این صورت که نمونه‌هایی که مقدایر Null دارند را به عنوان داده‌های تست و سایر نمونه‌ها را به عنوان داده‌های آموزش در نظر گرفته ایم. در گام دوم به بررسی ویژگی‌های موثر در پیش بینی بیماری قلبی پرداخته شد و با بهره وری از ۲ روش ضریب همبستگی پیرسون و PCA صفاتی که باهدف مسئله ارتباطی ندارند، شناسایی و حذف گردید. در گام ۳ با استفاده از ۳ الگوریتم SVM، J48 و بیزین پیش بینی باهدف تعیین بیماری قلبی بیمار ساخته شد که الگوریتم J48 با دقت ۸۹٪ دارای بالاترین دقت می‌باشد.

در کارهای آینده می‌توان از سایر روش‌های پیش بینی کننده استفاده کرد و به مقایسه با مدل پیشنهادی پرداخت. همچنین می‌توان با جمع آوری نمونه‌های بیشتر از مراکز مختلف بهداشتی فاکتورهای ریسکی که باعث شده افراد ساقه دار دوباره به این بیماری دچار شوند، را شناسایی نمود. موجود بودن اطلاعات دموگرافیک بیشتر و دقیق‌تر از بیماران از طریق مصاحبه حضوری از جمله سبک زندگی و تغذیه افراد و عامل مهم‌تر، قرارگیری فرد در شرایط استرس زا، نقش عوامل اجتماعی و روحی و روانی را در ایجاد بیماری قلبی-عروقی تا چه اندازه می‌تواند دخیل باشد، را بررسی و تحلیل نمود.

## فهرست و مراجع

- [۱] N. S. Chandra Reddy, S. Shue Nee, L. Zhi Min, and C. Xin Ying, "Classification and Feature Selection Approaches by Machine Learning Techniques: Heart Disease Prediction," *Int. J. Innov. Comput.*, vol. ۹, no. ۱, ۲۰۱۹.
- [۲] T. Karayilan and Ö. Kılıç, "Prediction of Heart disease using neural network," in *5nd International Conference on Computer Science and Engineering, UBMK ۲۰۱۷, ۲۰۱۷*, pp. ۷۱۹–۷۲۳.
- [۳] M. Raihan *et al.*, "Risk Prediction of Ischemic Heart Disease Using Artificial Neural Network," in *5nd International Conference on Electrical, Computer and Communication Engineering, ECCE ۲۰۱۹, ۲۰۱۹*, pp. ۱–۵.
- [۴] S. Radhimeenakshi and G. M. Nasira, "Prediction of Heart Disease using Neural Network with Back Propagation," *Int. J. Data Min. Tech. Appl.*, vol. ۴, no. ۱, pp. ۱۹–۲۲, ۲۰۱۵.
- [۵] S. K. Sen, "Predicting and Diagnosing of Heart Disease Using Machine Learning Algorithms," *Int. J. Eng. Comput. Sci.*, vol. ۶, no. ۶, Jun. ۲۰۱۷.
- [۶] S. Ansari *et al.*, "A review of automated methods for detection of myocardial ischemia and infarction using electrocardiogram and electronic health records," *IEEE Rev. Biomed. Eng.*, vol. ۱۰, pp. ۲۶۴–۲۹۸, ۲۰۱۷.
- [۷] H. Amin *et al.*, "Predictive Analysis of Heart Disease Using K-Means and Apriori Algorithms," *J. Appl. Sci. Comput.*, vol. VI, pp. ۲۱۸۳–۲۱۸۹, ۲۰۱۹.
- [۸] M. Mirmozaffari, A. Alinezhad, and A. Gilanpour, "Data Mining Apriori Algorithm for Heart Disease Prediction," *Int. J. Comput. Commun. Instrum. Eng.*, vol. ۴, no. ۱, ۲۰۱۷.
- [۹] P. C. M.A.JABBAR ,DEEKSHATULU, "Intelligent heart disease prediction system using random forest and evolutionary approach," *J. Netw. Innov. Comput.*, vol. ۴, pp. ۱۷۵–۱۸۴, ۲۰۱۶.
- [۱۰] A. H. Babolhakami, B. Behkamal, T. Dehghani, K. Etminani, and M. Naghibzadeh,

“Protein’s number of beta-sheets prediction using structural features,” in ۲۰۱۷ ۴th International Conference on Computer and Knowledge Engineering, ICCKE ۲۰۱۷, ۲۰۱۷, vol. ۲۰۱۷-Janua, pp. ۴۵۵–۴۶۰.

- [۱۱] T. Smayra, Z. Charara, G. Sleilaty, G. Boustany, L. Menassa-Moussa, and G. Halaby, “Classification and Regression Tree (CART) model of sonographic signs in predicting thyroid nodules malignancy,” *Eur. J. Radiol. Open*, vol. ۱, pp. ۳۴۳–۳۴۹, ۲۰۱۹.
- [۱۲] U. Sidiq, S. Mutahar Aaqib, and R. A. Khan, “Diagnosis of Various Thyroid Ailments using Data Mining Classification Techniques,” *Int. J. Sci. Res. Comput. Sci. Eng. Inf. Technol.*, pp. ۱۳۱–۱۳۶, ۲۰۱۹.
- [۱۳] UCI, “UCI Machine Learning Repository: Heart Disease Data Set,” *Uci*, ۲۰۱۹. [Online]. Available:  
<http://archive.ics.uci.edu/ml/datasets/Heart+Disease>. [Accessed: ۲۳-Nov-۲۰۱۹].