



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)
دانشکده مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی
گرایش مهندسی کامپیوتر

عنوان
پیش بینی و تشخیص بیماری های قلبی با استفاده از داده کاوی

نگارش
ارمغان سرور

استاد راهنما
دکتر رضا صفابخش

خردادماه ۱۳۹۸

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)
دانشکده مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی
گرایش مهندسی کامپیوتر

عنوان
پیش بینی و تشخیص بیماری های قلبی با استفاده از داده کاوی

نگارش
ارمغان سرور

استاد راهنما
دکتر رضا صفابخش

خردادماه ۱۳۹۸

سپاس‌گزاری

هرگاه که نگاه نویی نسبت به زندگی و عظمت آن می‌یابیم، میل به تلاش و امید در وجودمان زنده می‌شود.

وظیفه خود می‌دانم که مراتب امتنان خود را نسبت به استاد راهنمایم جناب آقای دکتر رضا صفا بخش که طی تدوین و پیش‌برد این پایان‌نامه مرا یاری نموده‌اند، اعلام بدارم.

ارمغان سرور
خرداد ۹۸

چکیده

امروزه کاربرد موفقیت‌آمیز داده‌کاوی در زمینه های بسیاری مانند تجارت الکترونیک، بازاریابی و صنایع دیگر قابل رؤیت است. در میان این بخش‌ها مراقبت‌های بهداشتی هنوز در مرحله کشف هستند. صنعت بهداشت و درمان به‌طور کلی شامل اطلاعات غنی است اما هنوز تمام اطلاعات برای کشف الگوهای پنهان و تصمیم‌گیری مؤثر استخراج نشده‌اند.

بیماری قلبی یکی از دلایل اصلی مرگ‌ومیر در سراسر جهان بوده و پیش‌بینی آن در مراحل اولیه ضروری است. سیستم‌های کامپیوتری می‌توانند به‌عنوان ابزار پیش‌بینی و تشخیص بیماری قلبی به پزشکان کمک کنند و با گسترش تحقیقات در این زمینه، آشنا نمودن خوانندگان با خلاصه‌ای از روش‌های فعلی پیش‌بینی و زیرمجموعه‌های آن‌ها اهمیت بالایی پیدا کرده است.

هدف این بررسی، شناسایی سیستم‌های پشتیبانی تصمیم برای پیش‌بینی و تشخیص این نوع بیماری‌ها بوده که از طریق داده‌کاوی و روش‌های هوشمند ترکیبی قابل بهبود است.

واژه‌های کلیدی:

داده‌کاوی، بیماری قلبی، سیستم پشتیبانی تصمیم

فهرست مطالب

۱	فصل اول: مقدمه.....
۲	مقدمه.....
۳	فصل دوم: تعریف و کاربرد داده کاوی.....
۴	تعریف و کاربرد داده کاوی.....
۴	۱-۲ تعریف داده کاوی.....
۵	۲-۲ تاریخچه.....
۶	۳-۲ اهمیت.....
۷	۲-۴ روش‌های داده کاوی.....
۸	۵-۲ کاربردها.....
۹	فصل سوم: درباره‌ی بیماری‌های قلبی.....
۱۰	درباره‌ی بیماری‌های قلبی.....
۱۰	۱-۳ تعریف.....
۱۱	۲-۳ وسعت بیماری قلبی.....
۱۳	فصل چهارم: الگوریتم‌های مورد استفاده.....
۱۴	الگوریتم‌های مورد استفاده.....
۱۴	۱-۴ درخت تصمیم.....
۱۶	۱-۱-۴ روش ID3.....
۱۶	۲-۱-۴ روش CART.....
۱۷	۳-۱-۴ روش J48.....
۱۹	۲-۴ الگوریتم Naïve Bayes.....
۲۰	۳-۴ اعمال الگوریتم ژنتیک.....
۲۱	فصل پنجم: جمع‌بندی و نتیجه‌گیری و پیشنهادها.....
۲۲	جمع‌بندی و نتیجه‌گیری و پیشنهادها.....
۲۲	۱-۵ جمع‌بندی.....
۲۳	۲-۵ نتیجه‌گیری.....
۲۴	۳-۵ پیشنهادها.....
۲۵	منابع و مراجع.....

فهرست اشکال

- فصل دوم: تعریف و کاربرد داده کاوی ۴
- شکل شماره ۱-۲ طبقه‌بندی الگوریتم‌های یادگیری ۷
- فصل چهارم: الگوریتم‌های مورد استفاده ۱۴
- شکل شماره ۱-۴ شمای کلی یک درخت تصمیم ۱۵
- شکل شماره ۲-۴ مراحل الگوریتم J48 ۱۸
- شکل شماره ۳-۴ اعمال الگوریتم Naïve Bayes به داده های جدید ۱۹
- شکل شماره ۴-۴ مراحل الگوریتم ژنتیک ۲۰

فهرست جداول و نمودارها

- فصل سوم: درباره بیماری‌های قلبی ۹
- جدول شماره ۱-۳ انواع بیماری‌های قلبی ۱۰
- نمودار شماره ۱-۳ نسبت مرگومیر ناشی از انواع بیماری‌های قلبی در ایران در سال ۱۳۹۵ ۱۲

فصل اول

مقدمه

مقدمه

قلب مهم‌ترین عضو عضلانی بدن انسان بوده که خون را به تمامی قسمت‌های آن پمپاژ می‌کند. زندگی انسانی وابسته به عملکرد مناسب قلب است و عملکرد نامنظم آن بر سایر بخش‌های بدن انسان نظیر مغز، کلیه و... تأثیر می‌گذارد. به‌طور کلی ایست جریان خون در قلب، حمله قلبی و ایست آن در مغز، سکته مغزی به‌حساب می‌آید. اگر گردش خون در بدن ناکارآمد باشد، هم بر عملکرد قلب و هم مغز تأثیر خواهد داشت. طبق گزارش آمار بهداشت جهانی علت اصلی ابتلا به بیماری و مرگومیر در جامعه مدرن بیماری قلبی است. برای مثال در کشور خودمان این رقم نزدیک به ۳۰ درصد بوده که نگرانی پزشکان و مهندسين حوزه بیوتکنولوژی را برمی‌انگیزد.

تشخیص پزشکی امری بسیار مهم اما پیچیده است که نیازمند تجربه، اطلاعات بالای پزشک و همچنین هزینه بالا از سوی بیمار است. اگرچه پیشرفت قابل‌توجهی در تشخیص و درمان بیماری‌های قلبی حاصل شده است اما تحقیقات باید به بالاترین صحت خود برسد.

دسترسی به تعداد زیادی از داده‌های پزشکی و گسترش داده‌های ذخیره‌شده منجر به نیاز به ابزارهای قدرتمند برای تجزیه و تحلیل آنها جهت استخراج دانش مفید است. داده‌کاوی از ابزارهای تحلیل مؤثر برای کشف روابط پنهان و الگوهای میان داده‌ها محسوب می‌شود. با بکارگیری الگوریتم‌های داده‌کاوی در حوزه تشخیص و درمان بیمارهای قلبی، می‌توان سیستم‌های هوشمندی ابداع نمود که به شکل خودکار و بدون نیاز به نظارت پزشک قادر به فهم و تفسیر بیمارهای قلبی افراد باشند و یا اطلاعات مفیدی را اکتشاف نمایند که متخصصان را در قضاوت صحیح یاری رساند.

در این گزارش تلاش گردیده تا کاربرد داده‌کاوی در این حوزه نشان داده شود و الگوریتم‌های سرآمد، معرفی و بررسی گردند.

فصل دوم

تعریف و کاربرد داده‌کاوی

تعریف و کاربرد داده کاوی

خوب است پیش از ورود به بحث اصلی، با داده کاوی بیشتر آشنا شویم. این شناخت به ما کمک می کند تا الگوریتم های داده کاوی را بهتر بشناسیم و برای بهبود آنها راه های بهتری پیشنهاد دهیم.

در ادامه بیشتر درباره ی داده کاوی خواهید خواند.

۱-۲ تعریف داده کاوی

به مجموعه ای از روش های قابل اعمال بر پایگاه داده های بزرگ و پیچیده به منظور کشف الگوهای پنهان و جالب توجه نهفته در میان داده ها، داده کاوی گفته می شود. روش های داده کاوی تقریباً همیشه به لحاظ محاسباتی پرهزینه هستند. علم میان رشته ای داده کاوی، پیرامون ابزارها، متدولوژی ها و تئوری هایی است که برای آشکارسازی الگوهای موجود در داده ها مورد استفاده قرار می گیرند و گامی اساسی در راستای کشف دانش محسوب می شود. دلایل گوناگونی پیرامون چرایی مبدل شدن داده کاوی به چنین حوزه مهمی از مطالعات وجود دارد. برخی از این موارد در ادامه بیان شده اند [۱].

- رشد انفجاری داده ها در گستره وسیعی از زمینه ها

- افزایش سریع قدرت پردازش کامپیوتری

روش های داده کاوی دارای انواع گوناگونی هستند و از کلاس بندی گرفته تا روش های تشخیص با الگوی پیچیده و هزینه محاسباتی بالا که ریشه در علوم کامپیوتر دارند را شامل می شوند. هدف اصلی این روش ها انجام پیش بینی می باشد اما این تنها هدف آن ها نیست.

۲-۲ تاریخچه

در سال ۱۹۶۰، کارشناسان آمار از اصطلاحات «صید داده»^۱ و «لایروبی داده»^۲ برای ارجاع به فعالیت‌های «تحلیل داده»^۳ استفاده می‌کردند. اصطلاح «داده‌کاوی» در حدود سال ۱۹۹۰ در جامعه پایگاه داده مورد استفاده قرار گرفت و به محبوبیت قابل‌توجهی دست پیدا کرد. عنوان مناسب‌تر برای فرآیند داده‌کاوی، «کشف دانش از داده»^۴ است.

در حال حاضر، یادگیری آماری و «علم داده»^۴ از دیگر عباراتی هستند که با معنای مشابه مورد استفاده قرار می‌گیرند، حال‌آنکه گاه تفاوت‌های ظریفی میان این موارد وجود دارد.

تولید داده موضوعی، به موازات پیشرفت‌های مادی و معنوی ابعاد گسترده‌ای یافته است و روز به روز نیز بر دامنه آن افزوده می‌شود. امروزه وسعت کاربرد آن کاملاً مشهود است و قرن ۲۱ را قرن تولید داده و فناوری اطلاعات می‌نامند.

^۱ Data Fishing

^۲ Data Dredging

^۳ Data Analytics

^۴ Data Science

۳-۲ اهمیت

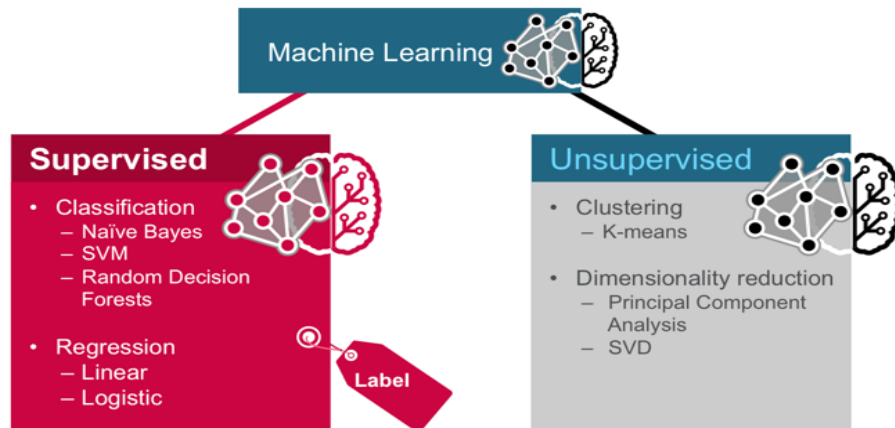
با رشد و افزایش توجه به داده‌کاوی، پرسش «چرا داده‌کاوی؟» مطرح می‌شود. در پاسخ به این پرسش باید گفت، داده‌کاوی کاربردهای زیادی دارد. بدین ترتیب، زمینه‌ای نوظهور و آینده‌دار برای نسل کنونی محسوب می‌شود. این زمینه توانسته توجهات زیادی را به صنایع و جوامع اطلاعاتی جلب کند که با وجود گستره وسیع داده‌ها، نیاز حتمی به تبدیل چنین داده‌هایی به اطلاعات و دانش وجود دارد.

بنابراین، بشر از اطلاعات برای دامنه‌ی وسیعی از کاربردها، از تحلیل بازار گرفته تا تشخیص بیماری‌ها، کشف کلاهبرداری و پیش‌بینی قیمت سهام استفاده می‌کند. در مجموع، باید گفت ضرب‌المثل انگلیسی «نیاز، مادر همه ابداعات بشر است»، پاسخی کوتاه و گویا به پرسش مطرح شده است.

۴-۲ روش‌های داده‌کاوی

در داده‌کاوی عمدتاً دو روش اصلی یادگیری وجود دارد. این دو رویکرد یادگیری نظارت‌شده^۱ و یادگیری بدون نظارت^۲ نام دارند که در ادامه به توضیح آن‌ها خواهیم پرداخت [۱].

- یادگیری نظارت‌شده: یادگیری تحت نظارت رایج‌ترین روش برای یادگیری است. در این روش از مدل‌های از پیش تعریف‌شده، با هدف آموزش استفاده می‌شود. مدل موردنظر با استفاده از داده‌های آموزش داده‌شده ساخته می‌شود. داده‌های ورودی جدید با مدل آموزش‌دیده مقایسه می‌شوند و برچسب کلاس داده‌های جدید تعیین خواهد شد. روش‌های طبقه‌بندی و رگرسیون زیرمجموعه‌ی این دسته هستند.
- یادگیری بدون نظارت: برخلاف روش‌های داده‌کاوی تحت نظارت، در این روش نتیجه‌ای از محیط اطراف دریافت نمی‌شود. اگرچه تجسم چگونگی آموزش یک دستگاه بدون پاسخ از محیط اطراف دشوار است اما این روش‌ها به‌خوبی کار می‌کنند. به احتمال زیاد، یک مدل مناسب برای روش‌های یادگیری بدون نظارت ایجاد می‌شود که هدف آن استفاده از ویژگی‌های ورودی فعلی برای پیش‌بینی ورودی‌های آینده، تصمیم‌گیری و انتقال مؤثر به یک الگوی دیگر است. خوشه‌بندی^۳ و روش کاهش ابعاد^۴ در این دسته قرار می‌گیرند.



شکل شماره ۱-۲ طبقه‌بندی الگوریتم‌های یادگیری

^۱ Supervised learning

^۲ Unsupervised learning

^۳ Clustering

^۴ Dimensionality reduction

۵-۲ کاربردها

داده‌کاوی در بسیاری از صنایع برای بهبود تجربه و رضایت مشتری و افزایش ایمنی و قابلیت استفاده محصول استفاده شده است. در مراقبت‌های بهداشتی نیز، استفاده از داده‌کاوی در زمینه‌هایی مانند پیش‌بینی بیماری، تشخیص تقلب و سوءاستفاده، مدیریت مراقبت‌های بهداشتی و اندازه‌گیری اثربخشی درمان‌های خاص اثبات شده است.

در ادامه به توضیح مختصری از دو مورد از این کاربردها می‌پردازیم.

- اندازه‌گیری اثربخشی درمان: این کاربرد از داده‌کاوی شامل مقایسه علائم، علل و دوره‌های درمان برای پیدا کردن مؤثرترین روش درمان یک بیماری خاص است. به‌عنوان مثال، گروه‌های بیمارانی که با رژیم‌های دارویی مختلف درمان می‌شوند، می‌توانند برای تعیین بهترین برنامه‌های درمان و صرفه‌جویی در هزینه‌ها، مقایسه شوند.
- تشخیص تقلب و سوءاستفاده: این مورد شامل ایجاد الگوهای طبیعی و سپس شناسایی الگوهای غیرمعمول ادعاهای پزشکی توسط کلینیک‌ها، پزشکان و یا آزمایشگاه‌ها است. این روش همچنین می‌تواند برای شناسایی ارجاع‌های نامناسب، نسخه‌های تقلبی و مدارک پزشکی جعلی مورد استفاده قرار گیرد.

فصل سوم

درباره‌ی بیماری‌های قلبی

درباره‌ی بیماری های قلبی

بیماری های قلبی و عروقی به معضلات زندگی انسان مدرن تبدیل شده اند. در نتیجه بخش مهمی از دانش بشری به این موضوع اختصاص یافته است. در این تحقیق نیز تمرکز ما بر بیماری قلبی است. لذا پیش از شروع بحث اصلی، ابتدا بهتر است با ابعاد و گسترش آن آشنا شویم.

۱-۳ تعریف

واژه بیماری قلبی معمولاً برای اشاره به حمله قلبی اشاره می شود، در حالی که شامل سایر مشکلات احتمالی قلب از جمله بیماری عروق کرونر، نارسایی قلبی و سکته قلبی می باشد. چند نمونه از انواع مختلف بیماری های قلبی همراه با شرح در جدول زیر آورده شده است.

جدول شماره ۱-۳: انواع بیماری های قلبی

توصیف	بیماری قلبی-عروقی
تأمین خون به عضله قلب به سرعت از بین می رود	سندرم حاد کرونری ^۱
درد قفسه سینه به علت کمبود خون ارسالی به عضله قلب	آنژین ^۲
شریان هایی که خون را به عضله قلب می رسانند، بسته می شوند	بیماری قلبی عروقی ^۳
رطوبت عضله یا بافت قلب	بیماری های التهاب قلب

مأخذ: سایت وزارت بهداشت و درمان، ۱۳۹۸

^۱ Acute coronary syndrome

^۲ Angina

^۳ Coronary heart disease

۲-۳ وسعت بیماری قلبی

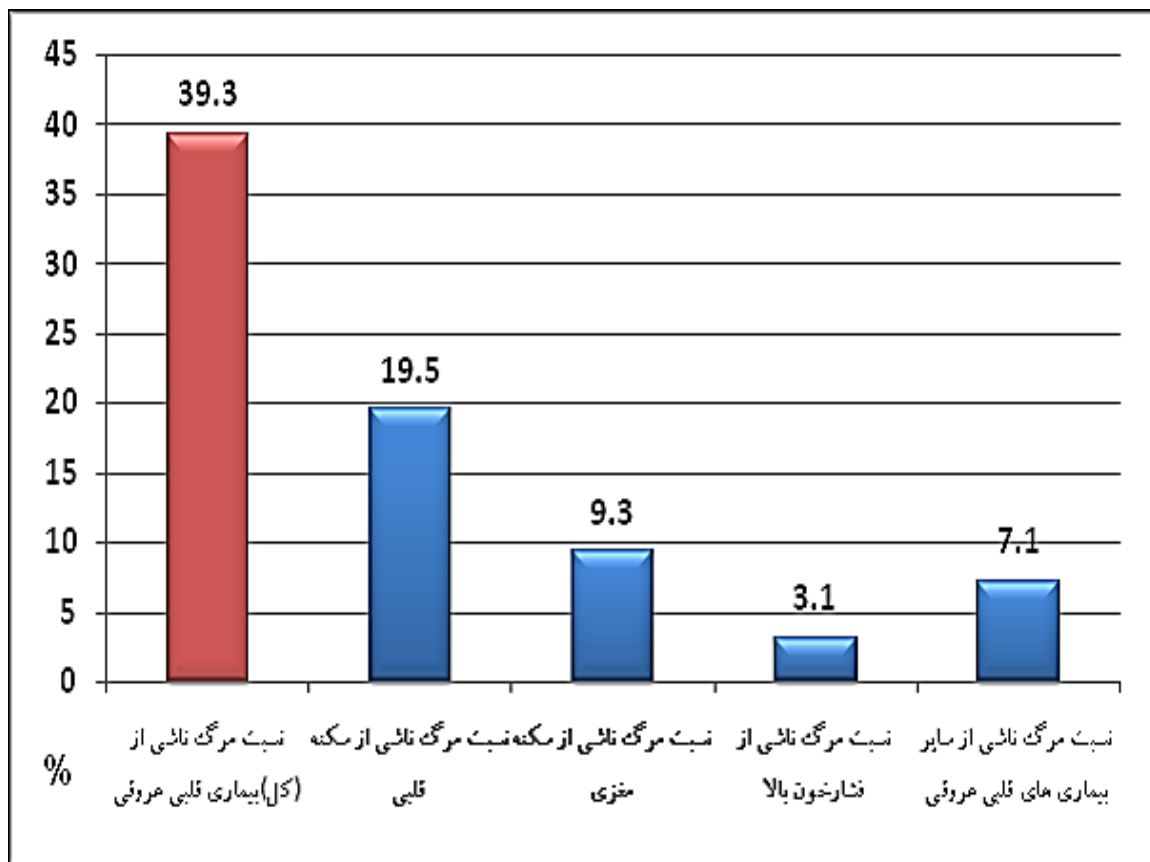
طبق برآورد صورت گرفته، ۱۷,۵ میلیون نفر در جهان در سال ۲۰۱۲ به علت بیماری های قلبی عروقی فوت نموده‌اند که ۳۱ درصد از کل موارد مرگومیرها را شامل می‌شود، از این مرگ‌ها حدود ۷,۴ میلیون به علت بیماری کرونر قلب و ۶,۷ میلیون ناشی از سکته‌ها بوده‌است. بیش از سه‌چهارم مرگومیر قلبی و عروقی در کشورهای کم‌درآمد و با درآمد متوسط رخ می‌دهد. از ۱۶ میلیون مرگومیر زیر ۷۰ سال به علت بیماری‌های غیر واگیر، ۸۲ درصد آن در کشورهای کم‌درآمد و با درآمد متوسط ایجاد می‌شود و ۳۷ درصد آن به علت بیماری‌های قلبی و عروقی است [۲].

متأسفانه آمار مرگومیر ناشی از بیماری‌های قلبی در کشور ما بالاتر از کشورهای پیشرفته است، بنابراین پیشگیری از این بیماری ضروری به نظر می‌رسد.

طبق آخرین آمارهای ارائه‌شده، ۵۰ درصد جامعه شهری ایران افزایش وزن دارند، این در حالی است که تحقیقات نشان می‌دهد چاقی شدید به‌تنهایی می‌تواند خطر نارسایی قلبی و مرگومیر ناشی از بیماری‌های قلبی را افزایش دهد.

بنابر آمارهای رسمی، مساله بیماری‌های قلبی یک مساله جهانی است و عدم تشخیص به موقع نارسایی در تشخیص، پیامدهای انسانی و مالی زیادی را متوجه بودجه‌های عمومی کشورها و خانوارها می‌سازد.

نمودار شماره ۳-۱ نسبت مرگومیر ناشی از انواع بیماری های قلبی در ایران در سال ۱۳۹۵



مأخذ: سایت وزارت بهداشت و درمان

فصل چهارم

الگوریتم‌های مورد استفاده

الگوریتم های مورد استفاده

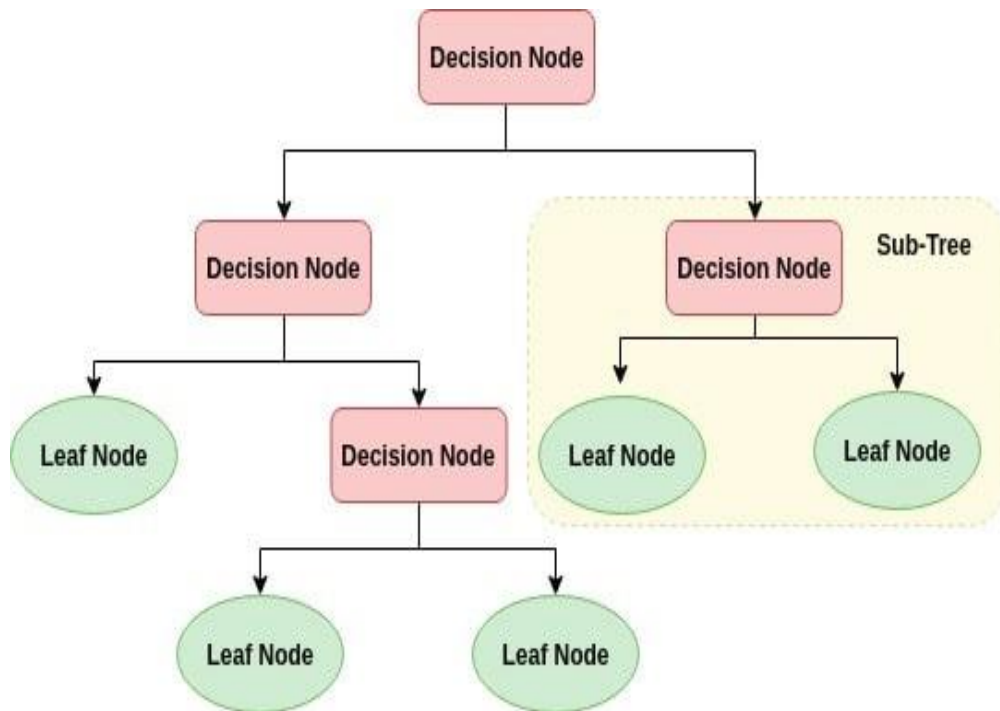
ابتدا نیاز داریم با روش های مختلف داده کاوی که تاکنون بر داده های بیماران اعمال گشته و آزمایش و بررسی شده اند آشنا شویم. اولین الگوریتم، درخت تصمیم^۱ می باشد که در ادامه به معرفی آن می پردازیم.

۴-۱ درخت تصمیم

درخت تصمیم یک طبقه بندی محبوب و ساده و آسان برای پیاده سازی محسوب می شود. در این روش هیچ نیازی به دانستن دامنه یا تنظیم پارامتر وجود ندارد و اجازه می دهد بتوانیم داده های با ابعاد بزرگ را مدیریت کنیم و نتایجی را که برای خواندن و تفسیر آسان هستند به دست آوریم. دسترسی به پروفایل دقیق بیماران تنها در درخت تصمیم امکان پذیر است.

یک درخت تصمیم گیری یک ساختار درخت مانند نمودار جریان است که هر گره داخلی نشان دهنده تست یک ویژگی و هر شاخه نشان دهنده نتیجه آزمایش است. همچنین گره های برگ نشان دهنده کلاس ها یا توزیع های کلاس هستند. گره بالاتر در یک درخت، گره ریشه است. گره های داخلی توسط مستطیل و گره های برگ توسط بیضی نشان داده می شوند. برای طبقه بندی یک نمونه ناشناخته، مقادیر ویژگی نمونه در درخت تصمیم گیری مورد آزمایش قرار می گیرند. مسیری از ریشه به گره برگ پیش بینی کلاس برای هر نمونه را نگه می دارد. بنابراین درخت های تصمیم می توانند به راحتی به قوانین طبقه بندی تبدیل شوند [۳].

¹ Decision Tree



شکل شماره ۴-۱ شمای کلی یک درخت تصمیم

انواع پیاده‌سازی‌های درخت تصمیم در موارد زیر متفاوت می‌شوند [۴]:

۱. معیار تقسیم (چگونگی محاسبه واریانس^۱)

۲. قابلیت ساخت مدل برای رگرسیون^۲

۳. نحوه مدیریت داده‌ها و اطلاعات ناقص

بر اساس این تفاوت‌ها الگوریتم‌های مختلفی شکل گرفتند که از میان آنها به موارد زیر اشاره می‌کنیم.

^۱ Variance

^۲ Regression

۴-۱-۱ روش ID3

این روش از اولین پیاده‌سازی‌های درخت تصمیم به شمار می‌رود که طی زمان اجرا، درخت تصمیم‌گیری برای طبقه‌بندی نمونه‌های جدید از پیش مشاهده نشده با کمک مقادیر گره‌های درخت مورد استفاده قرار گرفته و به شما کلاس توزیعی نمونه را برمی‌گرداند.

این الگوریتم جواب بهینه را تضمین نمی‌کند و ممکن است در بهینه‌های محلی به اتمام برسد اما به‌طور کلی درخت‌های کوچک تصمیم تولید می‌کند که استفاده از آن‌ها باعث عدم روی هم افتادگی کامل^۲ با داده‌ها خواهد شد [۵].

۴-۱-۲ روش CART

این روش، اگرچه عملکرد بخصوصی دارد اغلب به‌عنوان حالت کلی درخت تصمیم شناخته می‌شود. نمایش این مدل با استفاده از درخت دودویی^۴ انجام می‌شود. در این درخت هر گره ریشه بیانگر یک متغیر ورودی است و گره‌های برگ برای نمایش متغیر خروجی استفاده می‌شوند. در این الگوریتم ما از روابط زیر استفاده می‌کنیم:

• کسب اطلاعات^۵

$$\text{Information Gain} = I(p, n) \quad (۴-۱)$$

$$I(p, n) = \frac{-p}{p+n} \log_2 \left(\frac{p}{p+n} \right) - \left(\frac{n}{p+n} \right) \log_2 \left(\frac{n}{p+n} \right)$$

اما P و n در این رابطه بیانگر چه هستند؟ برای یافتن آن‌ها باید ویژگی یا نتیجه کلاس را مشخص کنیم که آن نیز دودویی است. برای p مقدار واقعی یک و برای n مقدار کاذب صفر را در نظر می‌گیریم.

¹ Iterative Dichotomiser 3

² Overfitting

³ Classification and Regression Tree

⁴ Binary

⁵ Information Gain

- بی‌نظمی^۱

$$Entropy(A) = \sum_{i=0}^v \frac{p_i + n_i}{p+n} (I(p, n)) \quad (۴-۲)$$

در حقیقت از بی‌نظمی برای ساخت درخت تصمیم استفاده می‌شود.

- بهره^۲

$$Gain = Information\ Gain - Entropy \quad (۴-۳)$$

$$Gain = I(p, n) - E(A)$$

و در نهایت، از بهره برای یافتن یکی از ویژگی‌های مجموعه آموزشی استفاده می‌کنیم.

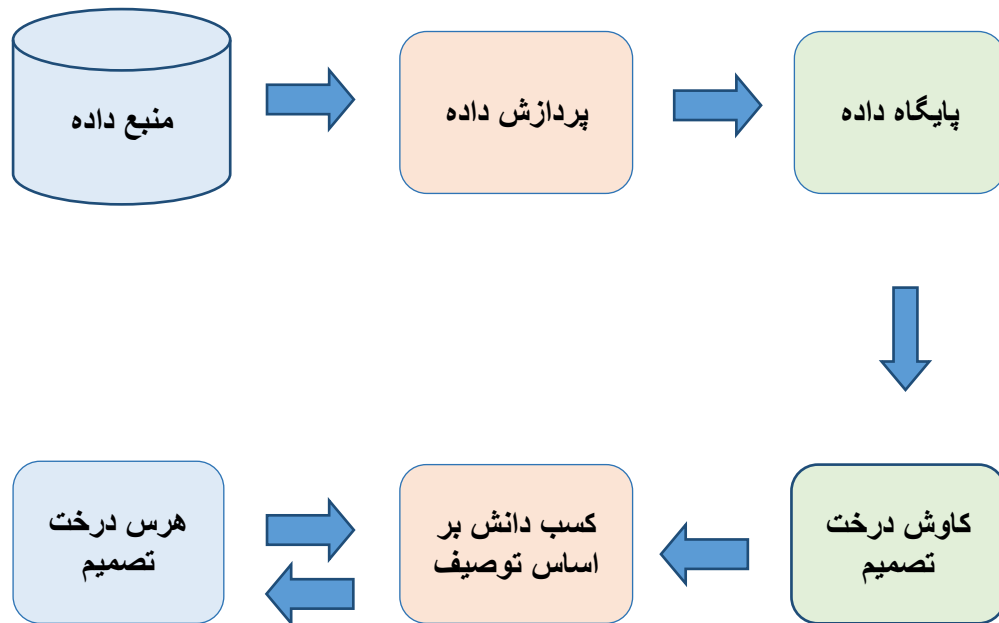
۴-۱-۳ روش J48

تصور کنید که شما یک مجموعه داده، لیست پیش‌بینی کننده‌ها یا متغیرهای مستقل و یک لیست از اهداف یا متغیرهای وابسته دارید. یک درخت تصمیم‌گیری مانند J48، به شما امکان می‌دهد که متغیر هدف یک نمونه داده جدید را پیش‌بینی کنید.

درواقع درخت تصمیم‌گیری J48 اجرای الگوریتم ID3 می‌باشد که توسط تیم پروژه WEKA توسعه یافته است. اجرای بازگشتی این الگوریتم و پردازش داده در مراحل اولیه، موجب تسریع در روند عملکرد آن شده است. در شکل زیر مراحل دقیق اجرا قابل مشاهده است [۶].

^۱ Entropy

^۲ Gain



شکل شماره ۴-۲ مراحل الگوریتم J48

۲-۴ الگوریتم Naïve Bayes

این الگوریتم یکی دیگر از سری الگوریتم های داده کاوی بوده که بر اساس اسم ابداع کننده اش نام گذاری شده است و الگوریتمی با نظارت محسوب می شود. این روش به طور کلی یک طبقه بندی آماری است که وابستگی بین صفات را در بر می گیرد و در عین حال از استقلال مشروط استفاده می کند. لذا هر داده در آن کاملاً مستقل از سایر داده ها مورد بررسی قرار می گیرد. در شکل زیر حالت کلی اعمال این روش بر داده ها قابل مشاهده است [۷].

Classification process

New data = $(X) = (X_1, X_2, \dots, X_m)$

Class C is a member of $\{C_1, C_2, \dots, C_k\}$

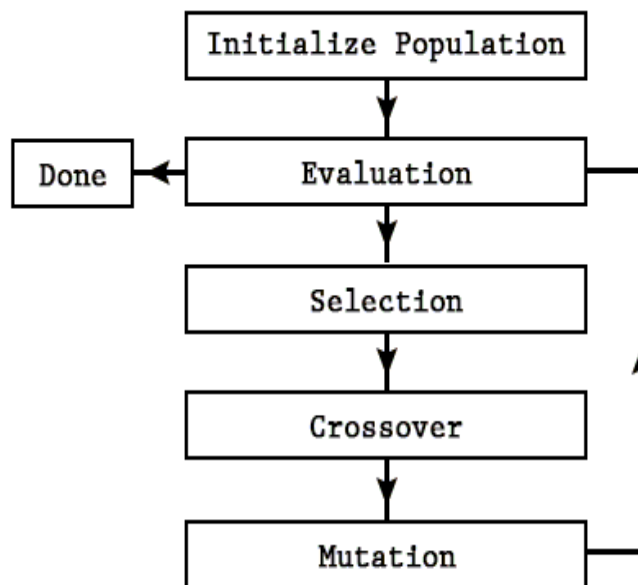


شکل شماره ۳-۴ اعمال الگوریتم Naïve Bayes به داده های جدید

۳-۴ اعمال الگوریتم ژنتیک

الگوریتم ژنتیک در هوش مصنوعی یک تکنیک جستجو است که از فرایند انتخاب طبیعی استفاده می‌کند. الگوریتم‌های تکاملی راحل‌هایی برای بهینه‌سازی و جستجوی مشکلات با استفاده از روش‌هایی مانند جهش، انتخاب و بازترکیبی ارائه می‌کنند. الگوریتم ژنتیک معمولی نیاز به نمایش ژنتیکی دامنه راحل و یک تابع تناسب^۱ برای ارزیابی دامنه راحل دارد. پس از انجام ارزیابی می‌توانیم اعضایی که شایستگی بالایی دارند انتخاب کرده و با اعمال ترکیب و جهش بر این اعضا تعادلی میان تکامل و تنوع که دو اصل مهم در الگوریتم‌های تکاملی محسوب می‌شوند برقرار کنیم [۸].

طبق آزمایش‌های اخیر، اگر الگوریتم‌های تکاملی مانند ژنتیک و یا الگوریتم‌های فازی را با سایر تکنیک‌های داده‌کاوی ترکیب کنیم، به جواب بهینه در مدت زمان کمتر و با سرعت بیشتری دست پیدا می‌کنیم.



شکل شماره ۴-۴ مراحل الگوریتم ژنتیک

^۱ Fitness function

فصل پنجم

جمع‌بندی و نتیجه‌گیری و پیشنهادها

جمع بندی و نتیجه گیری و پیشنهادها

۱-۵ جمع بندی

امروزه تحقیقات زیادی در زمینه‌ی داده‌کاوی در جریان است. ما نیز سعی کردیم در این تحقیق به بررسی چند مورد از این الگوریتم‌ها بپردازیم. هدف اصلی ساخت یک سیستم تشخیص خودکار پزشکی برای کمک به پزشکان و همچنین بیماران بود تا بتوانیم هزینه‌های درمان و همچنین خطای تشخیص بیماری را کاهش دهیم. از آنجا که بیماری‌های قلبی نیازمند تشخیص زودرس هستند، هدفمان ترکیب تکنیک‌های مطرح در حوزه پردازش داده و اعمال آن‌ها بر ویژگی‌های فیزیکی بیمار مانند وزن، جنس، فشارخون و... بود تا با جاسازی یک سیستم تشخیص خودکار در سازمان‌های بهداشتی و درمانی قادر باشیم پس از ورود یک بیمار و دریافت ویژگی‌های موردنظر، این داده‌ها را پردازش کرده و بیان کنیم بیمار با چه احتمالی به مشکلات قلبی دچار است. در نهایت با بررسی این الگوریتم‌ها کارایی آن‌ها را سنجیدیم و دیدیم که می‌توان با ترکیبشان با الگوریتم‌های فازی و تکاملی نتیجه را بهبود بخشید.

دیدیم که داده‌کاوی جستجوی خودکار منابع داده‌ای بزرگ، جهت یافتن الگوها و وابستگی‌هایی است که تحلیل‌های ساده آماری قادر به انجام آن‌ها نیستند. در علم پزشکی کشف و تشخیص به‌موقع بیماری‌ها می‌تواند از ابتلای افراد به بسیاری از بیماری‌های مهلک جلوگیری نموده و موجب نجات زندگی مردم گردد. این مطالعه نشان می‌دهد پیشگویی‌های داده‌کاوی ابزارهای ضروری را برای محققان، پزشکان و متخصصان حوزه بیماری‌های قلبی جهت بهبود در روش‌های تشخیصی و برنامه‌های درمانی فراهم می‌نمایند.

در سال‌های اخیر، جمع‌آوری داده‌های بیماری‌های مختلف اهمیت فراوانی یافته است. پیشرفت‌های مرتبط با فناوری اطلاعات کمک شایانی به بررسی‌های همه‌جانبه داده‌های حجیم به عمل آورده و توانسته است به جستجوی دانش نهفته در آن‌ها پرداخته و علم نوین داده‌کاوی را به وجود آورد.

۲-۵ نتیجه گیری

در جهان صنعتی امروز، سرعت به بخش لاینفک زندگی انسان ها تبدیل شده است. سرعت در تولید داده و شکل گیری بانک های عظیمی از آن ها، سبب شده است که روش ها و مدل های ساده دیگر کاربرد چندانی در تحلیل و استخراج دانش نداشته باشند. لذا داده کاوی برای جستجو، تحلیل و نتیجه گیری از بانک-داده های عظیم به روشی کارآمد و مؤثر تبدیل شده است.

یکی از زمینه هایی که نیازمند استفاده از ابزار های داده کاوی جهت مدل سازی پیش گویانه با روش های محاسباتی جدید است، علم پزشکی می باشد که کنترل و تشخیص بخش بزرگی از بیماری ها را شامل می شود. بیماری های قلبی به یکی از علل اصلی مرگومیر در جهان و ایران تبدیل شده اند و بخش بزرگی از داده های مؤثر از تشخیص، پیشگیری و درمان در این حوزه تولید می شود که الگوریتم های داده کاوی نقش محوری در مدیریت و تحلیل صحیح آنها خواهد داشت.

۳-۵-۳ پیشنهادها

در انتها پیشنهاد می‌شود درزمینه‌ی بهبود روش‌ها و الگوریتم‌های داده‌کاوی، نحوه استفاده از آن‌ها در سیستم پشتیبانی تصمیم و همچنین انواع دسته‌بندی ویژگی‌های بیماران مطالعه بیشتری صورت گیرد. این بهبود می‌تواند موجب افزایش دقت این سیستم شده و فرآیند تشخیص و پیش‌بینی بیماری را سرعت ببخشد.

منابع و مراجع

- [1] Tan steinbach and kumar Addison, 2nd ed. , “Introduction to Data Mining”, Wesley, 2006.
- [2] S. Vijayarani and S. Sudha, “A Study of Heart Disease Prediction in Data Mining”, *International Journal of Computer Science and Information Technology & Security (IJCSITS)* , vol. 2, no. 5, 2012.
- [3] Jyoti Soni, “Predictive Data Mining for Medical Diagnosis: An Overview of Heart Disease Prediction”, *International Journal of Computer Applications*, vol. 17, no. 8, 2011.
- [4] Chandna Deepali, “Diagnosis of heart disease using data mining algorithm” *International Journal of Computer Science and Information Technologies*, 5.2, 2015.
- [5] A. Wilson, G. Wilson and J. Likhiya, “Heart Disease Prediction using the Data Mining Techniques”, *International Journal of Computer Science Trends and Technology (IJCST)*, vol. 2, no. 1, 2014.
- [6] Vijayashree, J., and N. Ch SrimanNarayanaIyengar. “Heart disease prediction system using data mining and hybrid intelligent techniques: A review”, *International Journal of Bio-Science and Bio-Technology*, 2016.
- [7] Bhatla, Nidhi, and Kiran Jyoti, “An analysis of heart disease prediction using different data mining techniques”, *International Journal of Engineering*, 2012.
- [8] Shouman, M., T. Turner, and R. Stocker, “Using decision tree for diagnosing heart disease patients”, *9th Australasian Data Mining Conference*, 2011.