

بررسی عملکرد تکنیک‌های داده‌کاوی قواعد انجمنی و K همسایه نزدیک در پیش‌بینی میزان زنده ماندن افراد مبتلابه هپاتیت

محسن قربیان

گروه علوم کامپیوتر، واحد کاشان، دانشگاه آزاد اسلامی، کاشان، ایران.

نام نویسنده مسئول:

محسن قربیان

چکیده

گسترده‌ی تکنیک‌های داده‌کاوی موجود این امکان را فراهم کرده است تا بتوان با انتخاب تکنیک‌های مناسب، به نتایجی بهتر و با ضریب دقت بالاتر دست‌یافت؛ بنابراین با انتخاب تکنیک بهتر می‌توان به نتایج حاصل از داده‌کاوی صورت گرفته بر روی داده‌ها، اطمینان بیشتری حاصل نمود. وجود عوامل مختلف دخیل در بروز بیماری هپاتیت و همچنین ناقص بودن اطلاعات در دسترس را می‌توان از مشکلات اساسی در پیاده‌سازی تکنیک‌های داده‌کاوی دانست که می‌تواند بر روی نتایج نهایی تحقیق اثرگذار باشد. عوامل و نشانه‌های مختلفی در تشخیص بیماری هپاتیت به کار گرفته می‌شوند که از طریق انجام آزمایش‌های گوناگون این عوامل و نشانه‌ها مورد بررسی قرار می‌گیرند. در این تحقیق از بین عوامل که منجر به تشخیص بیماری هپاتیت می‌شوند نشانه‌هایی همانند ANOREXIA، LIVER، LIVER FIRM، BIG و SPLEEN PALPABLE مورد استفاده قرار می‌گیرند. از این رو با پیاده‌سازی تکنیک‌های داده‌کاوی قواعد انجمنی و K همسایه نزدیک با فاصله اطمینان ۸۰ درصد و مقدار پشتیبان ۵۰ درصد بر روی اطلاعات افراد مبتلابه هپاتیت و بر اساس معیارهای سنجش accuracy، error rate، specificity، negative prediction value عملکردشان مورد بررسی قرار گرفت و نتایج حاصل گویای این است که تکنیک قواعد انجمنی با مقدار 80/62 accuracy، 22/11 error rate، 95/27 specificity و negative prediction value 62/96 عملکرد بهتری را نسبت به تکنیک K همسایه نزدیک از خود نشان داده است. از بین عواملی که به‌عنوان نشانه‌های پیدایش بیماری هپاتیت محسوب می‌شوند عواملی که در نهایت فرد بیمار با دارا بودن آن‌ها می‌تواند شانس زنده ماندن را پیدا کند پیش‌بینی شدند. از این رو بر اساس نتایج به دست آمده می‌توان نتیجه گرفت که افرادی که در نتایج مرتبط با آزمایش خود با ویژگی‌های ANOREXIA، LIVER BIG، LIVER FIRM و SPLEEN PALPABLE مثبت مواجه شده است می‌توانند از شانس زنده ماندن برخوردار باشند.

واژگان کلیدی: داده‌کاوی پزشکی، قواعد انجمنی، K همسایه نزدیک، هپاتیت، تکنیک‌های داده‌کاوی.

مقدمه

به‌کارگیری کامپیوتر در علوم گوناگون و در ابعاد مختلف هم‌زمان با پیشرفت فناوری زمینه برای تولید بیش‌ازپیش اطلاعات دیجیتالی فراهم نموده است که این اطلاعات را می‌توان در قالب داده‌هایی با ساختار متناسب با خود در مخازنی تحت عنوان انبار داده ذخیره نمود [1]. ذخیره‌سازی داده‌ها در انبار داده‌ها را می‌توان فاز اول استفاده از داده‌های تولیدشده قلمداد کرد از این‌رو باید این اطلاعات ذخیره‌شده استفاده نمود. پردازش و استفاده از این اطلاعات با این حجم عظیم خارج از توان انسان است زیرا این اطلاعات دارای روابط و الگوهای نهانی است که تنها با استفاده از تکنولوژی‌ها و فناوری‌های ایجادشده به این منظور می‌توان آن‌ها را تحلیل کرد و مورد استفاده قرار داد. نتیجه بکارگیری این تحلیل‌ها را می‌توان برنامه‌ریزی برای آینده دانست که می‌تواند به نتایج مهمی منتهی گردد [2]. بنابراین می‌توان نتیجه گرفت که فناوری تحلیل اطلاعات از جمله فناوری‌های بسیار مهم و کاربردی است که می‌تواند به فرایند کشف دانش منتهی گردد از این‌رو فناوری بکار رفته در تحلیل اطلاعات را داده‌کاوی می‌نامند [3]. حوزه کاربرد داده‌کاوی^۱ را نمی‌توان به یک یا دو حوزه محدود کرد زیرا این فناوری را می‌توان در هر جایی که داده‌ای موجود باشد به کاربرد و از آن می‌توان به‌منظور آنالیز و تحلیل داده‌ها استفاده کرد در این‌بین داده‌هایی که مرتبط با حوزه پزشکی هستند هم این امر مستثنا نیستند و می‌توان آن‌ها را نیز مورد تحلیل و کاوش قرار داد [4]. برای تحلیل اطلاعات مرتبط با حوزه پزشکی باید به این امر دقت کرد که از آنجایی که این داده‌ها با سلامتی انسان‌ها در ارتباط هستند از این‌رو اگر تکنیک‌های داده‌کاوی بر روی این داده‌ها پیاده‌سازی شوند نتایج حاصل باید از دقت و صحت مناسبی برخوردار باشند [5]. با بکارگیری تکنیک‌های داده‌کاوی و پیاده‌سازی آن‌ها بر روی داده‌ها بیماران می‌توان به اطلاعات و الگوهای نهان موجود درون آن‌ها پی برد [6]. که نتایج کشف این الگوها می‌تواند در کاهش زمان تشخیص یک بیماری یا پیش‌بینی یک بیماری مورد استفاده قرار گیرد [7]. با توجه به گستردگی تکنیک‌های داده‌کاوی موجود لازم است تا تکنیک‌هایی را انتخاب کرد که بتوانند نتایجی با دقت و صحت مناسب را ارائه دهند از این‌رو باید با انتخاب تکنیکی که با تحلیل قوی خود قابلیت ارائه اطلاعاتی را داشته باشد که از صحت و دقت بالاتری برخوردار باشد [8]. برای به‌کارگیری فناوری داده‌کاوی نیز در حوزه سلامت نیاز است تا تکنیک‌هایی مورد استفاده قرار گیرند که بتوانند دقت و صحت نتایج حاصل را در بالاترین سطح ممکن به نمایش بگذارند تا بتوان به نتایج حاصل اطمینان کرد [9]. با توجه به اهمیتی که نتایج حاصله در پیاده تکنیک‌های داده‌کاوی بر روی داده‌های مرتبط به حوزه سلامت، بنابراین نیاز است تا تکنیک‌هایی که انتخاب می‌شوند از لحاظ ارائه نتایج قابل اطمینان باشند [10]. برای انتخاب تکنیکی که بتواند عملکرد مناسب و قابل اطمینانی از خود بروز دهد نیاز است تا این تکنیک‌ها مورد مقایسه قرار گیرند. مقایسه تکنیک‌ها از طریق معیارهایی صورت می‌گیرد که به معیارهای سنجش نیز معروف می‌باشند از این‌رو تکنیکی که بتواند با استفاده از این معیارهای سنجش نتایج قابل قبولی از خود به نمایش بگذارد را می‌توان به‌عنوان تکنیکی معرفی نمود که نسبت به تکنیک‌های دیگر در تحلیل اطلاعات و قابلیت اطمینان نتایج حاصل عملکرد بهتری از خود به نمایش گذاشته است [11]. هدف از این تحقیق بررسی عملکرد دو تکنیک داده‌کاوی K همسایه نزدیک^۲ و قواعد انجمنی^۳ بر اساس چهار معیار accuracy (معنای صحت)، error rate (نسبت خطا)، specificity (ویژگی) و negative prediction value (مقادیر منفی پیش‌بینی شده) است که می‌تواند منجر به انتخاب تکنیکی گردد که عملکرد مناسب‌تری را از خود به نمایش گذاشته است.

۱- پیشینه تحقیق

با گسترش علم و ظهور تکنولوژی‌های پیشرفته راه‌های حل مشکلات روزبه‌روز سهل و محمل تر شده است و این ویژگی منحصر به حوزه‌ی خاصی نیست و می‌توان ادعا کرد که استفاده از تکنولوژی در همه زمینه‌ها به‌عنوان یک اولویت در نظر گرفته می‌شود. حوزه پزشکی یکی از حوضه‌هایی است که انسان همواره با آن درگیر بوده است و متخصصین امر در این حوزه همیشه به دنبال افزایش کیفیت زندگی انسان‌ها بوده‌اند از این‌رو این متخصصین همیشه در تلاش بوده‌اند تا با مواردی همچون پیش‌گیری، افزایش سرعت تشخیص بیماری، کاهش هزینه‌های بیماری و کند کردن روند و سرعت رشد یک بیماری به هدف خود که همانا افزایش کیفیت زندگی انسان‌هاست نزدیک شوند. گسترش تکنولوژی در این حوزه نیز اثرات بخصوص خود را گذاشته است و باعث شده است تا با انجام آزمایش‌های گوناگون و یا به‌گونه‌ای دیگر با اختراع دستگاه‌های جراحی پیشرفته در خدمت بهبود کیفیت زندگی انسان باشند. فناوری داده‌کاوی را می‌توان یکی از تکنولوژی‌هایی دانست که با کاوش در اطلاعات مرتبط با بیماران و کشف رابطه میان این اطلاعات به الگوهای دست‌یابند که متخصصین امر در این حوزه یک دید استراتژیک بدهد تا بتوانند بر اساس روابط کشف‌شده میان این اطلاعات به الگوهای دست‌یابند که این الگوها می‌تواند آن‌ها را در جهت تشخیص سریع‌تر یا پیش‌بینی یک بیماری کمک کند که نتیجه آن می‌تواند به کاهش هزینه‌های

1. Data Mining
2. K-Nearest Neighbors
3. Rule Induction

تحمیلی به بیمار و همچنین افزایش اطمینان به روند بهبود سریع‌تر بیمار منجر گردد. طبیعی است که به دست آوردن اطلاعات با صحت و دقت بالا می‌تواند کمکی شایان به متخصصین نماید از این رو مهم است تا تکنیکی که انتخاب می‌گردد از نظر نتیجه حاصل، قابل اطمینان باشد. از این رو در سالیان اخیر متخصصین داده‌کاوی با همکاری متخصصین مربوط به حوزه پزشکی سعی کرده‌اند تا داده‌های مرتبط به بیماران را مورد تحلیل قرار دهند. در سال ۲۰۱۰ آقایان Rajeswari و Reena با استفاده از تکنیک‌های داده‌کاوی اقدام به تجزیه و تحلیل اطلاعات مرتبط با بیمارانی کردند که دچار مشکل کبدی بودند از این رو این دو متخصص با استفاده از نتایج حاصل از پیاده‌سازی تکنیک‌های داده‌کاوی بر روی ۳۴۵ نمونه توانستند میزان تأثیرگذاری مصرف الکل را در بروز اختلال در عملکرد کبد را پیش‌بینی کنند و سرعت رشد بیماری را در بیمارانی که الکل مصرف کرده‌اند را پیش‌بینی کنند [12]. اما در سال ۲۰۰۷ آقای Nguyen به همراه همکاران توانستند که الگوریتم داده‌کاوی جدیدی را معرفی کنند که این الگوریتم با استفاده از الگوهای ایجادشده بین زمانی که افراد در آن دوره زمانی به بیماری هپاتیت مبتلا شده‌اند و همچنین عوامل مؤثر در بروز بیماری توانسته‌اند تا به الگویی دست یابند که این الگو به پزشکان کمک خواهد کرد تا بتوانند زمان تشخیص این بیماری را بهبود بخشند [13]. در سال ۲۰۰۷ آقای Kim و همکاران با استفاده از تکنیک‌های داده‌کاوی ماشین بردار و درخت تصمیم‌گیری توانستند تا میزان حساسیت کبدهایی که درگیر هستند و امکان تبدیل شدن این حساسیت به بیماری مزمن هپاتیت را مورد بررسی قرار دادند و سعی کردن تا عملکرد این دو تکنیک را در تشخیص این مهم مورد ارزیابی قرار دهند [14]. در سال ۲۰۱۳ آقای Radwan و همکاران با استفاده از تکنیک داده‌کاوی درخت تصمیم و پیاده‌سازی آن بر روی داده‌های بیماران مبتلا به هپاتیت توانستند مدل‌هایی را پیش‌بینی کنند که به واسطه آن مدل‌ها توانستند نتیجه استفاده از داروهای ضد ویروسی مربوط به هپاتیت را پیش‌بینی کنند که این تحقیق می‌تواند منجر به کاهش هزینه‌های تحمیلی بر بیمار گردد [15]. در سال ۲۰۰۵ آقای Yokoi و همکاران با پیاده‌سازی تکنیک‌های داده‌کاوی بر روی اطلاعات به دست آمده از آزمایش‌های ادرار بیماران مبتلا به هپاتیت و تحلیل نتایج حاصل توانستند به مدل و الگویی دست یابند که می‌تواند به متخصصان امر در حوزه درمان هپاتیت کمک شایانی نماید و روند درمان این بیماری را تسریع بخشد [16]. در سال ۲۰۰۲ آقای Sato و همکاران با پیاده‌سازی تکنیک درخت تصمیم بر روی اطلاعات بیماران مبتلا به هپاتیت و تحلیل نتایج به دست آمده توانستند به الگوهایی دست یابند. سپس این محققان با ترکیب الگوهای به دست آمده توانستند به مکانیزمی دست یابند که می‌توانست کار متخصصان را بسیار آسان نماید زیرا متخصصان برای به دست آوردن چنین نتایجی از طریق آزمایش، نیاز به انجام آزمایش‌های بسیار سخت و طولانی مدت داشتند [17].

۲- روش پژوهش

یکی از چالش‌های اساسی برای محققان در فرایند داده‌کاوی انتخاب تکنیکی است که متناسب با حوزه فعالیت مربوطه باشد. به عنوان نمونه هنگامی داده‌های مورد استفاده مربوط به حوزه پزشکی و سلامت است اهمیت نتایج حاصل دوچندان خواهد شد زیرا بر اساس این اطلاعات متخصصین امر در مورد نحوه درمان یا تشخیص یک بیماری تصمیم‌گیری می‌کنند که این مهم می‌تواند در سلامت فرد تأثیر مستقیمی بگذارد. از این رو نیاز است تا تکنیک‌هایی انتخاب شوند که بتوانند بالاترین سطح از اطمینان و درستی را در نتایج به دست آمده ارائه کنند تا بتوان درست‌ترین تصمیم را با بهترین تحلیل صورت گرفته را همراه نمود و نتیجه را به بهترین شکل ممکن دریافت کرد. در این تحقیق سعی شده است تا دو تکنیک K همسایه نزدیک و قواعد انجمنی را به عنوان تکنیک‌های منتخب برگزید که دلیل این امر را می‌توان گستردگی استفاده از این دو تکنیک و همچنین مطلوب بوده کیفیت تحلیل‌های ارائه شده از جانب این دو تکنیک دانست. از مدل ایجادشده توسط این تکنیک به منظور نمایش حقایق استنتاج شده مربوط به متغیر هدف استفاده می‌شود. مدل‌های ارائه شده توسط قواعد انجمنی بسیار ساده و قابل فهم می‌باشند و در قالب قوانینی ساده بیان می‌شوند که از این رو تفسیر این قوانین به سادگی صورت می‌گیرد [18]. تکنیک K همسایه نزدیک داده‌ها را در گروه‌هایی تقسیم‌بندی می‌کند که دارای ویژگی‌های نزدیک به یکدیگر می‌باشند و در صورتی که نمونه‌ی جدیدی را بخواهد مورد بررسی قرار دهد این نمونه را در قسمت گروه‌هایی قرار می‌دهد که به نمونه مورد نظر بیشترین شباهت را داشته باشد. این تکنیک را می‌توان به عنوان یکی از ساده‌ترین تکنیک‌ها در مبحث طبقه‌بندی داده‌ها در نظر گرفت. تکنیک K همسایه نزدیک سعی می‌کند با قرار دادن نمونه‌ای جدید و با استفاده از اطلاعات طبقه‌بندی داده‌های قبلی، با قرار دادن نمونه‌های جدید در طبقه‌بندی نزدیک به آن‌ها سعی در افزایش دقت و تحلیل‌های حاصل نماید [19]. هپاتیت^۴ در لغت به معنی متورم شدن است که این امر می‌تواند دلایل گوناگونی را داشته باشد، مصرف دخانیات، مصرف مشروبات الکلی و یا مواد شیمیایی می‌توانند از دلایل متورم شدن کبد محسوب شوند. هپاتیت دارای انواع گوناگونی است که از این انواع می‌توان به موارد زیر اشاره نمود:

هپاتیت نوع A

هپاتیت نوع B

هپاتیت نوع C

تفاوت انواع هپاتیت را می‌توان در نوع ویروسی که باعث التهاب و متورم شدن کبد می‌شود بیان کرد. بیماری‌های هپاتیت A و هپاتیت B دارای واکسن می‌باشند که انسان می‌تواند با واکسینه کردن خود از ابتلا به این بیماری در امان باشد ولی از طرف دیگر بیماری هپاتیت C فاقد واکسن است و نمی‌توان از طریق واکسینه کردن از ابتلا به آن در امان بود. از این رو این امکان نیز وجود داد که یک فرد به‌طور هم‌زمان به دو یا سه نوع هپاتیت مبتلا گردد. بیماری هپاتیت C را می‌توان از خطرناک‌ترین انواع بیماری هپاتیت محسوب کرد. راه انتقال این نوع بیماری از طریق خون به‌صورت مستقیم و یا خال‌کوبی، مصرف مواد مخدر و طب سنتی قابل‌انتقال است. [20]. داده‌های مورد استفاده در این تحقیق مربوط به ۱۵۵ بیمار مبتلا به هپاتیت است. این اطلاعات توسط دانشگاه کالیفرنیا در ارواین تهیه شده‌اند که به‌عنوان یکی از معتبرترین دانشگاه‌های ایلات متحده آمریکا محسوب می‌شود. اطلاعات موجود در این داده‌ها را می‌توان به قسمت‌های مختلفی دسته‌بندی کرد. این اطلاعات در مورد افرادی می‌باشند که بر روی آن‌ها آزمایش‌های یکسانی صورت گرفته‌اند و در انتها ضمن مشخص کردن جنسیت آن‌ها میزان مرگ‌ومیر افراد مربوطه را به نمایش گذاشته است. در این تحقیق سعی بر آن است تا با استفاده از تکنیک‌های داده‌کاوی ضمن مقایسه عملکرد تکنیک‌های مطرح‌شده بر اساس معیارهای سنجش تعریف‌شده، با بررسی و تحلیل نتایج حاصل احتمال زنده ماندن فردی که آزمایش‌های مربوطه بر روی آن صورت گرفته است پیش‌بینی شود. در این بین به علت تکرار آزمایش‌های صورت گرفته سعی بر آن شده است تا چهار آزمایش را به‌عنوان متغیرهای این تحقیق در نظر بگیریم و این چهار متغیر تحت لقای کلاس زنده ماندن یا فوت کردن بیمار قرار می‌گیرند. این چهار متغیر را می‌توان به این صورت نمایش داد:

- ANOREXIA (بی اشتهاپی)

- LIVER BIG (بزرگی کبد)

- LIVER FIRM (کبد چرب)

- SPLEEN PALPABLE (قابل لمس بودن طحال)

روش انجام تحقیق در (شکل ۱) نشان داده شده است.



شکل ۱: مکانیزم پیاده‌سازی تحقیق

۳-۱- پردازش اولیه

به‌منظور استفاده درست و بهینه از تکنیک‌های داده‌کاوی نیاز است تا ساختار اطلاعات وارد شده به این تکنیک‌ها استاندارد و مطابق با ساختار در نظر گرفته شده برای آن‌ها باشد. از این داده‌های مرتبط با بیماران هپاتیت نیز از این قاعده مستثنا نیستند و می‌بایست دارای ساختاری متناسب با تکنیک‌های دسته‌بندی استفاده شده در داده‌کاوی باشند. پس نیاز است تا در اطلاعات خام اولیه تغییراتی ایجاد شود.

۳-۱-۱- تبدیل اطلاعات بیماران به قالب استاندارد

تکنیک‌های داده‌کاوی در نظر گرفته شده برای این تحقیق از دسته تکنیک‌های رده‌بندی^۵ می‌باشند از این‌رو نیاز است تا اطلاعات بیماران هپاتیت ساختارمند گردند و آن‌ها را در قالب سطر و ستون نمایش داد. از طرفی به دلیل این‌که این اطلاعات دارای ساختار ابتدایی بوده‌اند پس نیاز است با اندکی تغییر آن‌ها را به فرم مورد نظر تبدیل نمود.

۳-۱-۲- تبدیل نتایج حاصل از آزمایش به Y و N

بعد از ساختاردهی ابتدایی به اطلاعات نیاز است تا در مقادیر موجود تغییراتی اعمال شود. اطلاعات بیماران مربوط به آزمایش‌های صورت گرفته بر روی آن‌ها است، از این رو نتایج به‌صورت اعداد نمایش داده شده‌اند و از طرفی این اعداد دارای مقادیر ۱ و ۲ می‌باشند از این‌رو با توجه به توضیحات مندرج در ارتباط با اطلاعات این بیماران نتایج آزمایش‌ها به‌صورت Yes و No در نظر گرفته شده‌اند و ۱ به این

معنا است که در صورتی که در مقابل متغیری مقدار ۱ قرار گرفته باشد یعنی این که جواب آزمایش فرد برای آن متغیر مقدار No بوده است و در صورتی که در مقابل متغیری مقدار ۲ قرار گرفته باشد یعنی جواب آن آزمایش Yes بوده است. از این رو به منظور استفاده بهینه و افزایش عملکرد تکنیک‌های داده‌کاوی و خوانایی و درک ساده تر نتایج حاصله به جای مقادیر ۲ و ۱ به ترتیب مقادیر Y و N که بیانگر Yes و No می‌باشند را قرار می‌دهیم.

۲-۲-۳- انتخاب ویژگی‌ها و آماده‌سازی داده‌ها

اطلاعات موجود شامل آزمایش‌های گوناگونی است. با توجه به کثرت آزمایش‌های صورت گرفته و گستردگی بازه‌های تعریف شده در برخی متغیرها از این رو نیاز است تا متغیرهای مورد نظر را از میان سایر متغیرهای انتخاب کنیم تا بتوانیم تکنیک‌های داده‌کاوی را تنها بر روی متغیرهای انتخاب شده پیاده‌سازی کنیم.

۲-۲-۳-۱- انتخاب ویژگی‌ها

تعداد آزمایش‌های صورت گرفته بر روی بیماران مبتلا به هیپاتیت ۲۰ مورد است. طبیعی است بررسی این ۲۰ مورد آزمایش خارج از هدف تعیین شده برای این تحقیق است و ما تنها ۴ مورد را به عنوان متغیر انتخاب می‌کنیم و تکنیک‌های داده‌کاوی را بر روی آن‌ها پیاده‌سازی می‌کنیم. دلیل انتخاب این چهار متغیر را می‌توان اهمیت و فراوانی آن‌ها و همچنین تأثیرگذاری بیشتر این آزمایش‌ها در روند درمان بیمار در نظر گرفت. از این رو ویژگی‌های انتخاب شده را می‌توان در جدول ۱ مشاهده کرد.

جدول ۱: انتخاب ویژگی‌ها

Value		Properties
LIVE	DIE	CLASS
P	N	ANOREXIA
P	N	LIVER BIG
P	N	LIVER FIRM
P	N	SPLEEN PALPABLE

۲-۲-۳-۲- جایگذاری مقادیر ناموجود

از بین ویژگی‌های انتخاب شده به عنوان متغیر، فیلد^۶ مربوط به بعضی از آزمایش‌ها فاقد مقادیر می‌باشند از این رو عدم وجود این مقادیر می‌تواند بر روی نتایج و تحلیل‌های حاصل از پیاده‌سازی تکنیک‌های داده‌کاوی تأثیرگذار باشند. از این رو باید برای این فیلدهای خالی مقادیری را در نظر گرفت تا بتوان تأثیر ناشی از عدم وجود بعضی از مقادیر را از بین برد. به این منظور از امکانی که داده‌کاوی در قالب نرم‌افزار پیاده ساز آن یعنی Rapid Miner در اختیار ما قرار داده است استفاده می‌کنیم. امکان در نظر گرفته شده راه‌های گوناگونی را به منظور پر کردن مقادیر خالی فیلدها به ما پیشنهاد می‌دهد ولی در این تحقیق از روش میانگین‌گیری استفاده می‌شود به این صورت که از تمامی مقادیر موجود در ستون مربوط به ویژگی انتخاب شده میانگین گرفته می‌شود و نتیجه حاصل در خانه‌های خالی از مقدار مربوطه قرار می‌گیرد. در (جدول ۲) می‌توان مقادیر ناموجود در هر یک از ویژگی‌های انتخاب شده را مشاهده کرد.

جدول ۲: مقادیر ناموجود ویژگی‌ها انتخاب شده

Null Value	Properties
1	ANOREXIA
10	LIVER BIG
11	LIVER FIRM
5	SPLEEN PALPABLE

۳-۳- مدل‌سازی

مدل‌سازی از داده‌ها را می‌توان به عنوان بخشی معرفی نمود که بعد از مرحله آماده‌سازی داده اتفاق خواهد افتاد و در این مرحله تکنیک‌های انتخاب شده بر روی داده‌هایی با ساختار مناسب و استاندارد متناسب پیاده‌سازی خواهند شد. فرایند مدل‌سازی که تکنیک‌های داده‌کاوی آن را انجام می‌دهند به این صورت خواهد بود که در ابتدای امر این تکنیک‌ها نیاز دارند که به یک دانش قبلی دست یابند و با

استفاده از این دانش قبلی خود اقدام به مدل‌سازی کنند یعنی به‌عبارت‌دیگر بتوانند فرایند پیش‌بینی را انجام دهند. در این قسمت داده‌کاوی امکانی را فراهم کرده است تا به‌واسطه آن ما بتوانیم برای تکنیک‌های داده‌کاوی مشخص کنیم که از چه مقدار از داده‌های ورودی به‌عنوان داده‌های یادگیری استفاده کند و این یعنی همان دانش قبلی که تکنیک‌های داده‌کاوی به آن نیاز دارند و در ادامه تکنیک‌های داده‌کاوی بادانشی که از این قسمت به دست آورده‌اند اقدام به انجام پیش‌بینی بر روی مابقی اطلاعات موجود خواهند نمود. این امر واضح و مبرهن است که هر مقداری دانش تکنیک‌های داده‌کاوی در مورد اطلاعات بالاتر باشد امکان پیش‌بینی بهتر و دقیق‌تری مهیا می‌گردد.

۳-۳-۱- انتخاب تکنیک‌های داده‌کاوی

با توجه به گستردگی تکنیک‌های داده‌کاوی موجود از رو نیاز است تا تکنیک‌های موردنظر خود را انتخاب نمایم تا بتوانیم آن‌ها را بر روی داده‌های خود پیاده‌سازی کنیم. از این‌رو نیاز است تا داده‌های خود را متناسب با تکنیک‌های در نظر گرفته‌شده آماده‌سازی کنیم و بر اساس آن تکنیک‌ها ساختار قالب‌ها را تعریف نمایم. تکنیک‌های استفاده‌شده در این تحقیق از نو تکنیک‌های دسته‌بندی‌شده‌اند به‌عبارت‌دیگر داده‌های موردنیاز برای پیاده‌سازی این دسته از تکنیک‌ها نیاز است که دارای سطر و ستون باشند.

۳-۳-۲- ساخت مدل

بعد از آماده‌سازی داده‌ها و انتخاب تکنیک‌های متناسب حال نیاز است تا فرایند مدل‌سازی صورت پذیرد. از این مرحله می‌توان به‌عنوان آخرین مرحله از داده‌کاوی یاد کرد. در این مرحله با پیاده‌سازی تکنیک‌های داده‌کاوی بر روی داده‌های موجود می‌توان به مدل‌های ایجادشده توسط تکنیک‌های داده‌کاوی دست‌یافت. از این‌رو بعد از نتایج به‌دست‌آمده لازم است تا این نتایج را در قالب معیارهای سنجشی که از قبل در نظر گرفته‌شده‌اند مورد مقایسه قرارداد تا به این شکل بتوان تکنیکی را که عملکرد بهتری را از خود نشان داده است را معرفی نمود. معیارهای تعریف‌شده در این تحقیق عبارت‌اند از accuracy, specificity, error rate, negative prediction value که با استفاده از این چهار معیار سنجش عملکرد دو تکنیک داده‌کاوی K همسایه نزدیک و قواعد انجمنی را می‌تواند ارزیابی کرد.

۳-۳-۲-۱- accuracy

accuracy یا به‌عبارت‌دیگر معیار صحت را می‌توان به‌عنوان یکی از مهم‌ترین و پراستفاده‌ترین معیارهای سنجش یک مدل معرفی کرد. معیار صحت از طریق محاسبه مجموع تعداد مواردی که کلاس آن‌ها توسط مدل به‌درستی پیش‌بینی شده‌اند به تعداد کل مواردی که توسط مدل پیش‌بینی شده‌اند به دست می‌آید. فرمول محاسبه معیار صحت از قرار زیر است [21].

$$accuracy = \frac{TP+TN}{TP+FN+FP+TN} \quad (1)$$

۳-۳-۲-۲- error rate

error rate یا معیار نسبت خطا معکوس عملکرد معیار صحت رفتار می‌کند. معیار نسبت خطا از طریق محاسبه مجموع تعداد مواردی که کلاس آن‌ها توسط مدل به‌اشتباه پیش‌بینی شده‌اند به تعداد کل مواردی که توسط مدل پیش‌بینی شده‌اند به دست می‌آید. به‌عبارت‌دیگر اگر نتیجه حاصل از معیار صحت را از مقدار عددی ۱ کم کنیم مقدار معیار نسبت خطا به دست خواهد آمد. فرمول محاسبه معیار نسبت خطا از قرار زیر است [22].

$$error\ rate = \frac{FN+FP}{TP+FN+FP+TN} \quad (2)$$

۳-۳-۲-۳- specificity

specificity یا ویژگی را به‌عنوان یکی از معیار مهم در سنجش عملکرد تکنیک‌های داده‌کاوی مورد استفاده قرار می‌گیرد. معیار ویژگی از طریق محاسبه نسبت تعداد موارد منفی حقیقی تخمین زده‌شده توسط مدل به مجموع موارد مثبت کاذب و منفی حقیقی پیش‌بینی شده به دست می‌آید. فرمول محاسبه معیار ویژگی به‌قرار زیر تعریف می‌گردد. [23]

$$specificity = \frac{TN}{FP+TN} \quad (3)$$

۳-۳-۲-۴- negative prediction value

negative prediction value که به‌اختصار NPV نامیده می‌شود معیار سنجشی است که به نسبت زیاد مورد استفاده قرار می‌گیرد. NPV از طریق محاسبه نسبت موارد منفی حقیقی پیش‌بینی شده به مجموع موارد منفی کاذب و منفی حقیقی پیش‌بینی شده به دست می‌آید. فرمول NPV به‌قرار زیر است. [24]

$$NPV = \frac{TN}{FN+TN} \quad (4)$$

اجزای تشکیل‌دهنده رابطه‌های ۴ تا ۱ هر یک به‌گونه‌ای بیان‌کننده حالاتی از پیش‌بینی‌های صورت گرفته می‌باشند که از طریق اعمال تکنیک‌ها بر روی نمونه‌ها حاصل شده‌اند. این اجزا را می‌توان به‌صورت زیر تعریف کرد.

TP: تعداد نمونه‌هایی که به‌درستی مثبت تشخیص داده می‌شوند.

TN: تعداد نمونه‌هایی که به‌درستی منفی تشخیص داده می‌شوند.

FP: تعداد نمونه‌هایی که به‌اشتباه مثبت تشخیص داده می‌شوند.

FN: تعداد نمونه‌هایی که به‌اشتباه منفی تشخیص داده می‌شوند.

۳-۵- نتیجه نهایی

هدف از پیاده‌سازی این تحقیق بررسی عملکرد دو تکنیک داده‌کاوی K همسایه نزدیک و قواعد انجمنی است که این مقایسه و ارزیابی با استفاده از چهار معیار سنجش accuracy, error rate, specificity, negative prediction value صورت می‌پذیرد؛ و در ادامه تکنیکی که دارای عملکرد مناسب‌تری باشد انتخاب می‌شود و از طرف دیگر با تحلیل نتایج حاصل از پیاده‌سازی تکنیک‌های داده‌کاوی بر روی داده‌های بیماران مبتلابه هپاتیت سعی می‌شود تا عواملی که در صورت مثبت شدن در آزمایش فرد مبتلابه هپاتیت، فرد می‌تواند امید زنده ماندن را داشته باشد، مورد بررسی قرار می‌گیرند.

۴- نتایج پژوهش

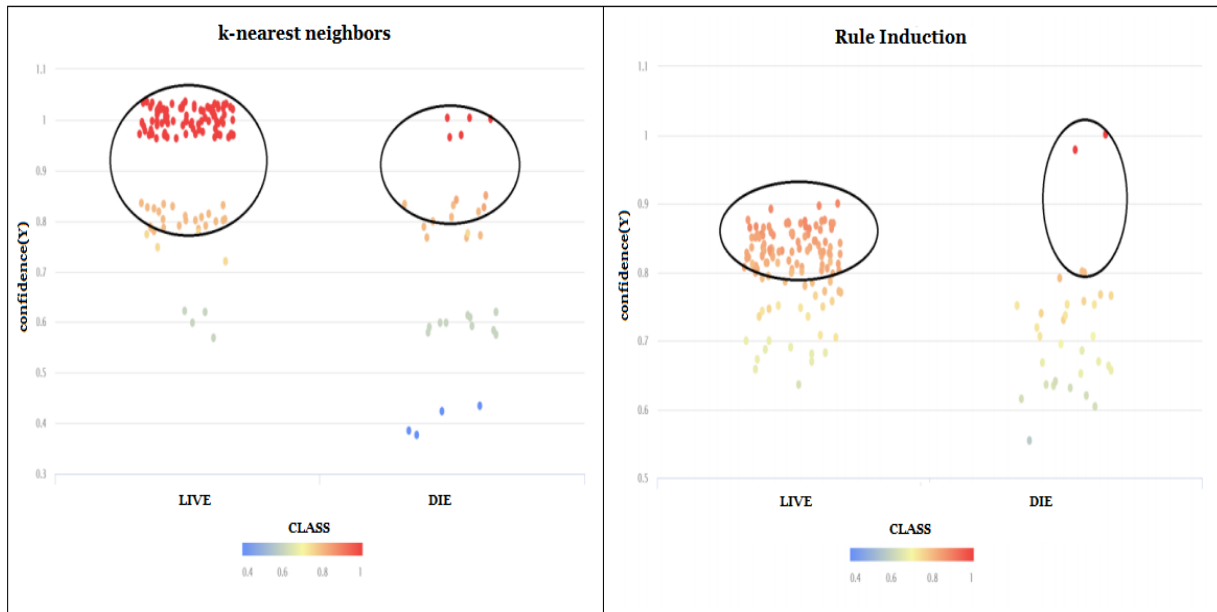
در این بخش در ابتدای امر سعی شده است تا نتایج حاصل از پیاده‌سازی دو تکنیک K همسایه نزدیک و قواعد انجمنی بر روی چهار ویژگی ANOREXIA, LIVER BIG, LIVER FIRM و SPLEEN PALPABLE مورد بررسی قرار گیرد. در این بخش دو قسمت جواب وجود خواهد داشت قسمت اول شامل نتایج عملکرد تکنیک‌های داده‌کاوی پیاده‌سازی شده K همسایه نزدیک و قواعد انجمنی خواهد بود که بر اساس چهار معیار accuracy, error rate, specificity, negative prediction value مورد ارزیابی قرار می‌گیرند و در این قسمت تکنیکی که عملکرد مطلوب‌تری را که از خود نشان داده‌شده باشد معرفی می‌گردد و در قسمت دوم مربوط به نتیجه سعی شده است تا تأثیر مثبت بودن هر یک از چهار ویژگی تعریف‌شده ANOREXIA, LIVER BIG, LIVER FIRM و SPLEEN PALPABLE را در زنده ماندن فرد مبتلابه هپاتیت مورد بررسی قرار گیرد. در این قسمت به‌منظور تأثیرگذاری یک ویژگی شروطی لحاظ شده است تا در صورتی که یک ویژگی بتواند آن شروط را احراز کند به‌عنوان یک ویژگی در نظر گرفته می‌شود که فرد مبتلابه هپاتیت در صورتی که از آن برخوردار باشد می‌تواند امید زنده ماندن داشته باشد به‌عبارت‌دیگر در صورتی که نتیجه مربوط به آن آزمایش مثبت باشد فرد مبتلابه همچنان می‌تواند امید به حیات داشته باشد. ویژگی‌های تعریف‌شده برای این تحقیق مقدار فاصله اطمینان ۸۰ درصد و مقدار پشتیبان ۵۰ درصد است. علت انتخاب مقدار ۸۰ درصد برای فاصله اطمینان را می‌توان به دلیل حساسیت بالای حوزه پزشکی در نظر گرفت و با توجه به مطرح‌شدن موضوع با سلامت انسان نیاز است تا مقدار فاصله اطمینان بالا در نظر گرفته شود تا نتایج حاصل‌شده از دقت و اطمینان کافی برخوردار باشند. مقدار پشتیبان ۵۰ درصد به این معنی است که در صورتی که از ۱۵۵ نمونه موجود ویژگی در حداقل ۷۸ مورد مثبت تشخیص داده‌شده است و در این حالت فرد نیز زنده مانده باشد در این صورت است که این ویژگی به‌عنوان مشخصه‌ای در نظر گرفته می‌شود که فرد در صورتی که در آزمایش‌های صورت گرفته مقداری مثبت برای این ویژگی به‌دست‌آمده باشد فرد نیز می‌تواند همچنان امید به حیات داشته باشد.

جدول ۳: وضعیت متغیرهای انتخاب شده در ۱۵۵ نمونه موجود.

قواعد انجمنی				K همسایه نزدیک		DIE				LIVE		۱۵۵ نمونه موجود			پارامتر بالینی
NPV	ERR	SPEC	ACC	NPV	ERR	SPEC	ACC	مثبت	منفی	مثبت	منفی	مثبت	منفی	مقادیر ناموجود	
(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	تعداد	تعداد	تعداد	تعداد	تعداد	تعداد		
۷۸,۸	۲۳,۲۱	۹۶,۸	۷۶,۷	۷۹,۳	۲۰,۶۷	100	۷۹,۳	۱۹	۱۰	۱۰۳	2۲	۲12	32	۱	ANOREXIA
۲۰	۱۸,۰۸	۹۳,۸	۸۱,۹	۱۲,۵	۱۹,۹۶	۹۰,۳۳	۸۰,۴	۲۲	۳	۹۸	۲۲	۱۲۰	25	۱۰	LIVER BIG
۷۲,۴	۲۷,۷۹	۹۰,۴۴	۷۲,۲	۷۳,۵	۲۶,۵۸	۹۰,۴۴	۷۳,۴	۱۴	۱۱	۷۲	۴۷	۸۶	۵۸	۱۱	LIVER FIRM
۸۰,۶	۱۹,۳۷	۱۰۰	۷۷,۸	۸۰,۱	۱۹,۳۷	۱۰۰	۷۸,۸	۱۹	۱۲	۱۰۲	۱۰	۱۲۱	۲۹	۵	SPLEEN PALPABLE

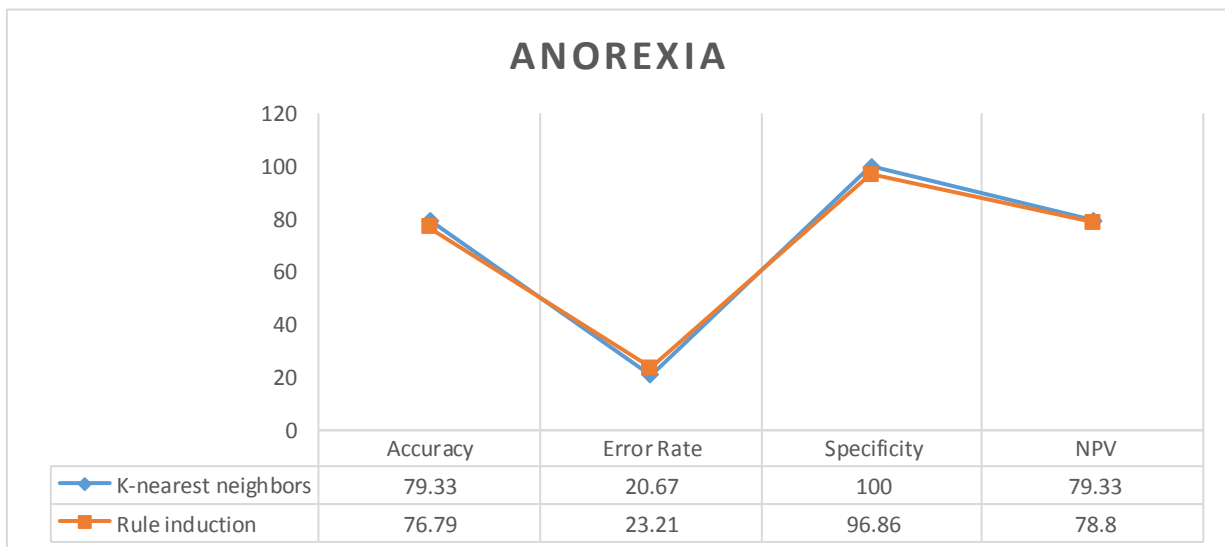
۴-۱- ویژگی ANOREXIA

با استناد به نتایج به دست آمده در صورتی که فرد مبتلا به بیماری هپاتیت در نتیجه آزمایش‌ها خود با ANOREXIA مثبت مواجه شود پیش‌بینی می‌شود فرد بتواند زنده بماند (شکل ۲).



شکل ۲: نمودار عملکرد تکنیک‌های داده‌کاوی در پیش‌بینی زنده ماندن بیمار در صورت مثبت شدن ANOREXIA

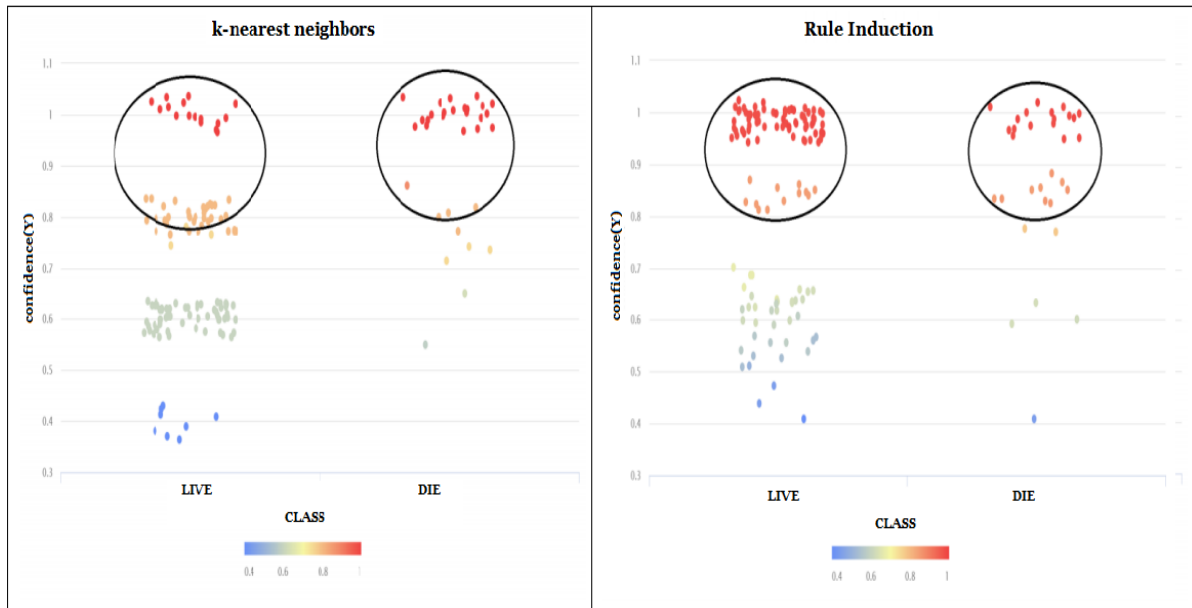
با استناد به نتایج به‌دست‌آمده از پیاده‌سازی تکنیک‌های داده‌کاوی و بر اساس ارزیابی‌های صورت گرفته با استفاده از چهار معیار سنجش accuracy، error rate، specificity و negative prediction value می‌توان بیان داشت که تکنیک داده‌کاوی K همسایه نزدیک نسبت به تکنیک قواعد انجمنی در پیش‌بینی زنده ماندن فرد مبتلابه هیپاتیت، در صورتی‌که فرد نتیجه ANOREXIA مثبت را در آزمایش خود داشته باشد، دارای عملکرد بهتری است (شکل ۳).



شکل ۳: ارزیابی عملکرد تکنیک‌های داده‌کاوی در پیش‌بینی زنده ماندن بیمار در صورت مثبت شدن ANOREXIA

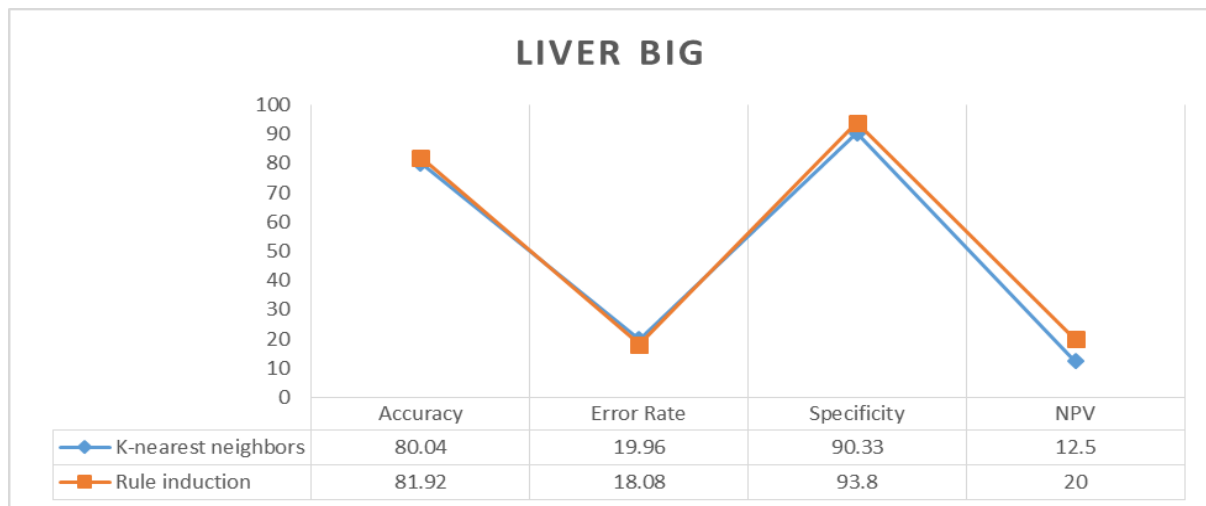
۲-۴- ویژگی LIVER BIG

با استناد به نتایج به‌دست‌آمده در صورتی‌که فرد مبتلابه بیماری هیپاتیت در نتیجه آزمایش‌ها خود با LIVER BIG مثبت مواجه شود پیش‌بینی می‌شود فرد بتواند زنده بماند (شکل ۴).



شکل ۴: نمودار عملکرد تکنیک‌های داده‌کاوی در پیش‌بینی زنده ماندن بیمار در صورت مثبت شدن LIVER BIG

با استناد به نتایج به‌دست‌آمده از پیاده‌سازی تکنیک‌های داده‌کاوی و بر اساس ارزیابی‌های صورت گرفته با استفاده از چهار معیار سنجش accuracy، error rate، specificity و negative prediction value می‌توان بیان داشت که تکنیک داده‌کاوی قواعد انجمنی نسبت به تکنیک K همسایه نزدیک در پیش‌بینی زنده ماندن فرد مبتلابه هیپاتیت، در صورتی که فرد نتیجه LIVER BIG مثبت را در آزمایش خود داشته باشد، دارای عملکرد بهتری است (شکل ۵).



شکل ۵: ارزیابی عملکرد تکنیک‌های داده‌کاوی در پیش‌بینی زنده ماندن بیمار در صورت مثبت شدن LIVER BIG

۳-۴- ویژگی LIVER FIRM

با استناد به نتایج به‌دست‌آمده در صورتی که فرد مبتلابه بیماری هیپاتیت در نتیجه آزمایش‌ها خود با LIVER FIRM مثبت مواجه شود پیش‌بینی می‌شود فرد بتواند زنده بماند (شکل ۶).



شکل ۶: نمودار عملکرد تکنیک‌های داده‌کاوی در پیش‌بینی زنده ماندن بیمار در صورت مثبت شدن LIVER FIRM

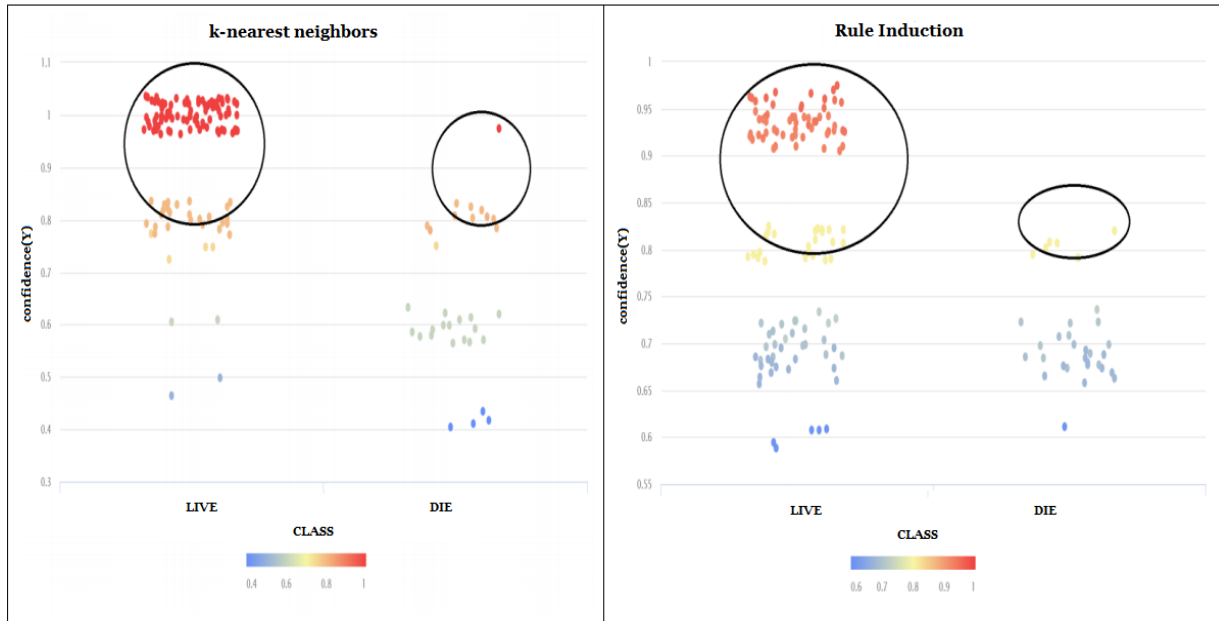
با استناد به نتایج به‌دست‌آمده از پیاده‌سازی تکنیک‌های داده‌کاوی و بر اساس ارزیابی‌های صورت گرفته با استفاده از چهار معیار سنجش accuracy، error rate، specificity و negative prediction value می‌توان بیان داشت که تکنیک داده‌کاوی K همسایه نزدیک نسبت به تکنیک قواعد انجمنی در پیش‌بینی زنده ماندن فرد مبتلا به هپاتیت، در صورتی که فرد نتیجه LIVER FIRM مثبت را در آزمایش خود داشته باشد، دارای عملکرد بهتری است (شکل ۷).



شکل ۷: ارزیابی عملکرد تکنیک‌های داده‌کاوی در پیش‌بینی زنده ماندن بیمار در صورت مثبت شدن LIVER FIRM

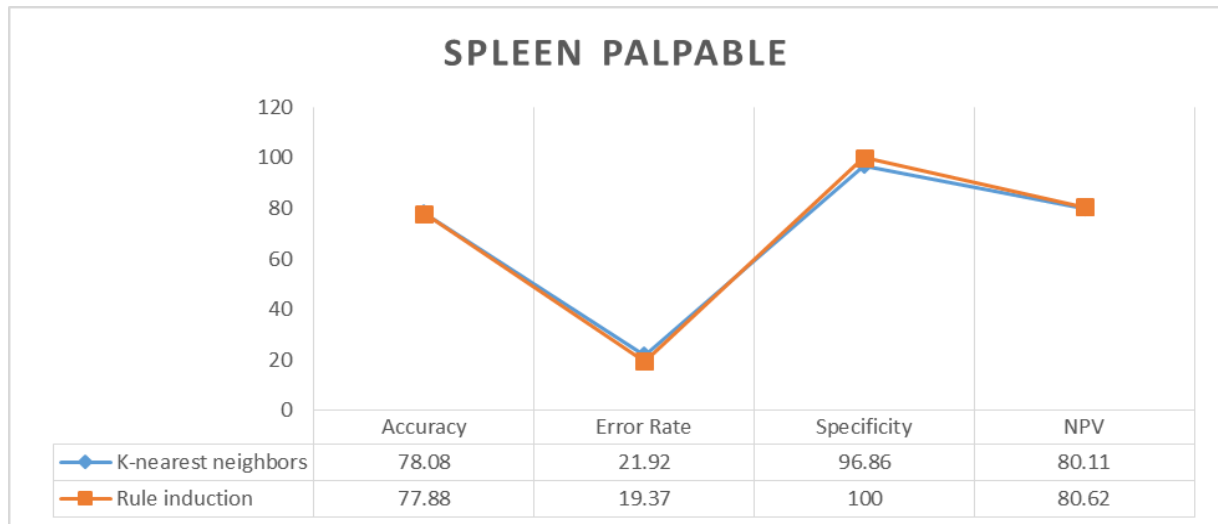
۴-۴- ویژگی SPLEEN PALPABLE

با استناد به نتایج به‌دست‌آمده در صورتی که فرد مبتلا به بیماری هپاتیت در نتیجه آزمایش‌ها خود با SPLEEN PALPABLE مثبت مواجه شود پیش‌بینی می‌شود فرد بتواند زنده بماند (شکل ۸).



شکل ۸: نمودار عملکرد تکنیک‌های داده‌کاوی در پیش‌بینی زنده ماندن بیمار در صورت مثبت شدن SPLEEN PALPABLE

با استناد به نتایج به‌دست‌آمده از پیاده‌سازی تکنیک‌های داده‌کاوی و بر اساس ارزیابی‌های صورت گرفته با استفاده از چهار معیار سنجش accuracy, error rate, specificity و negative prediction value می‌توان بیان داشت که تکنیک داده‌کاوی قواعد انجمنی نسبت به تکنیک K همسایه نزدیک در پیش‌بینی زنده ماندن فرد مبتلا به هیپاتیت، در صورتی که فرد نتیجه LIVER BIG مثبت را در آزمایش خود داشته باشد، دارای عملکرد بهتری است (شکل ۹).



شکل ۹: ارزیابی عملکرد تکنیک‌های داده‌کاوی در پیش‌بینی زنده ماندن بیمار در صورت مثبت شدن SPLEEN PALPABLE

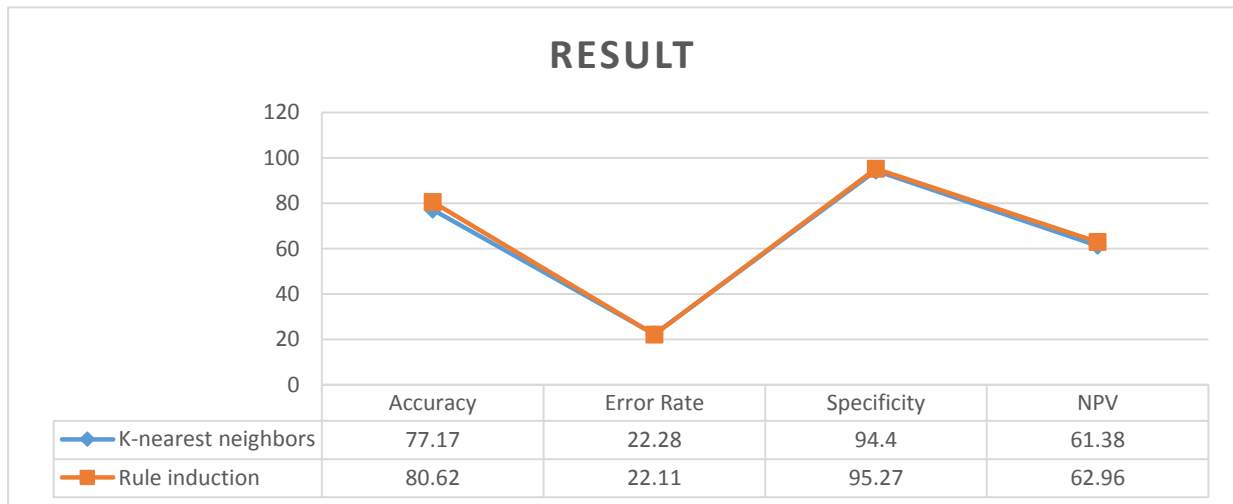
۵- نتیجه‌گیری نهایی

با جمع‌آوری نتایج حاصل از پیاده‌سازی دو تکنیک داده‌کاوی قواعد انجمنی و K همسایه نزدیک بر روی چهار ویژگی LIVER FIRM, LIVER BIG, ANOREXIA و SPLEEN PALPABLE که در (جدول ۳) به نمایش درآمده‌اند، اقدام به میانگین گرفتن از این نتایج می‌شود. با میانگین گرفتن از نتایج حاصل می‌توان به عملکرد کلی تکنیک‌ها در طول فرآیندهای انجام‌شده دست‌یافت. نتایج حاصل از میانگین‌گیری در (جدول ۴) نشان داده شده‌اند.

جدول ۴: مقایسه نهایی عملکرد دو تکنیک داده‌کاوی

تکنیک‌های داده‌کاوی	accuracy	rate error	specificity	NPV
همسایه نزدیک K	77.17	22.28	94.4	61.38
قواعد انجمنی	80.62	22.11	95.27	62.96

نتایج به‌دست‌آمده نشان می‌دهد قواعد انجمنی در مجموع ارزیابی‌های صورت گرفته دارای عملکرد بهتری در مقایسه با K همسایه نزدیک است (شکل ۱۰).



شکل ۱۰: ارزیابی نهایی از عملکرد دو تکنیک داده‌کاوی

در مجموع پیاده‌سازی‌های صورت گرفته با استفاده از تکنیک‌های داده‌کاوی K همسایه نزدیک و قواعد انجمنی بر روی داده‌های مرتبط با بیماران مبتلابه هپاتیت و بررسی عملکرد این دو تکنیک بر اساس چهار معیار ارزیابی accuracy, error rate, specificity و negative prediction value می‌توان نتیجه گرفت که تکنیک قواعد انجمنی با مقدار accuracy 80/62, error rate 22/11, specificity 95/27 و negative prediction value 62/96 عملکرد بهتری را نسبت به تکنیک K همسایه نزدیک از خود نشان داده است؛ و در مورد نتایج حاصل از بررسی‌های صورت گرفته بر روی ویژگی‌های موجود در آزمایش افراد مبتلابه هپاتیت می‌توان نتیجه گرفت که افرادی که در نتایج مرتبط با آزمایش خود با ویژگی‌های ANOREXIA, LIVER BIG, LIVER FIRM و SPLEEN PALPABLE مثبت مواجه شده است می‌توانند از شانس زنده ماندن برخوردار باشند.

سپاسگزاری

بدینوسیله نویسنده از داوران ارجمند برای نظرات ارزنده و پیشنهادهای در راستای بهبود ارائه این مقاله، تشکر میکند.

منابع و مراجع

- [1] Ng RT, Pei J. Introduction to the special issue on data mining for health informatics. *ACM SIGKDD Explorations Newsletter*. 2007; 9(1):1-2.
- [2] Soni J, Ansari U, Sharma D, Soni S. Predictive data mining for medical diagnosis: An overview of heart disease prediction. *International Journal of Computer Applications*. 2011; 17(8):43-8.
- [3] Wu X, Kumar V, Quinlan JR, Ghosh J, Yang Q, Motoda H, et al. Top 10 algorithms in data mining. *Knowledge and information systems*. 2008; 14(1):1-37.
- [4] Padhy N, Mishra D, Panigrahi R. The survey of data mining applications and feature scope. *ArXiv preprint arXiv: 12115723*. 2012.
- [5] Wu X, Zhu X, Wu G-Q, Ding W. Data mining with big data. *IEEE transactions on knowledge and data engineering*. 2014; 26(1):97-107.
- [6] Delen D, Walker G, Kadam A. Predicting breast cancer survivability: a comparison of three data mining methods. *Artificial intelligence in medicine*. 2005; 34(2):113-27.
- [7] Glover S, Rivers PA, Asoh DA, Piper CN, Murph K. Data mining for health executive decision support: an imperative with a daunting future! *Health services management research*. 2010; 23(1):42-6.
- [8] Gharehchopogh FS, Molany M, Mokri FD. Using artificial neural network in diagnosis of thyroid disease: a case study. *International Journal on Computational Sciences & Applications (IJCSA) Vol.* 2013; 3:49-61.
- [9] Fayyad UM, Piatetsky-Shapiro G, Smyth P, Uthurusamy R. *Advances in knowledge discovery and data mining*. 1996.
- [10] Maroco J, Silva D, Rodrigues A, Guerreiro M, Santana I, de Mendonça A. Data mining methods in the prediction of Dementia: A real-data comparison of the accuracy, sensitivity and specificity of linear discriminant analysis, logistic regression, neural networks, support vector machines, classification trees and random forests. *BMC research notes*. 2011; 4(1):299.
- [11] Sokolova M, Japkowicz N, Szpakowicz S, editors. *Beyond accuracy, F-score and ROC: a family of discriminant measures for performance evaluation*. Australasian joint conference on artificial intelligence; 2006: Springer.
- [12] Rajeswari P, Reena GS. Analysis of liver disorder using data mining algorithm. *Global journal of computer science and technology*. 2010.
- [13] Ho T-B, Nguyen C-H, Kawasaki S, Le S-Q, Takabayashi K. Exploiting temporal relations in mining hepatitis data. *New Generation Computing*. 2007;25(3):247.
- [14] Uhm S, Kim D-H, Cho SW, Cheong JY, Kim J, editors. *Chronic hepatitis classification using SNP data and data mining techniques*. 2007 *Frontiers in the Convergence of Bioscience and Information Technologies*; 2007: IEEE.
- [15] Zayed N, Awad AB, El-Akel W, Doss W, Awad T, Radwan A, et al. The assessment of data mining for the prediction of therapeutic outcome in 3719 Egyptian patients with chronic hepatitis C. *Clinics and research in hepatology and gastroenterology*. 2013;37(3):254-61.
- [16] Abe H, Ohsaki M, Yokoi H, Yamaguchi T, editors. *Implementing an integrated time-series data mining environment based on temporal pattern extraction methods: a case study of an interferon therapy risk mining for chronic hepatitis*. Annual Conference of the Japanese Society for Artificial Intelligence; 2005: Springer.
- [17] Ohsaki M, Sato Y, Yokoi H, Yamaguchi T, editors. *A rule discovery support system for sequential medical data, in the case study of a chronic hepatitis dataset*. Workshop Notes of the International Workshop on Active Mining, at IEEE International Conference on Data Mining; 2002.
- [18] Yin X, Han J, editors. *CPAR: Classification based on predictive association rules*. Proceedings of the 2003 SIAM International Conference on Data Mining; 2003: SIAM.

- [19] Guo G, Wang H, Bell D, Bi Y, Greer K, editors. KNN model-based approach in classification. OTM Confederated International Conferences" On the Move to Meaningful Internet Systems"; 2003: Springer.
- [20] Lauer GM, Walker BD. Hepatitis C virus infection. New England journal of medicine. 2001;345(1):41-52.
- [21] Lucas P. Bayesian analysis, pattern analysis, and data mining in health care. Current opinion in critical care. 2004; 10(5):399-403.
- [22] Witt O, Deubzer HE, Milde T, Oehme I. HDAC family: What are the cancer relevant targets? Cancer letters. 2009; 277(1):8-21.
- [23] Cios KJ, Moore GW. Uniqueness of medical data mining. Artificial intelligence in medicine. 2002;26(1-2):1-24.
- [24] Steinberg DM, Fine J, Chappell R. Sample size for positive and negative predictive value in diagnostic research using case-control designs. Biostatistics. 2008;10(1):94-105.